



# Toward scalable solvers for generalized n-body problems

## **Georgia Tech**

– *Aparna Chandramowlishwaran, Ilya Lashuk, Ryan Riegel, Aashay Shringarpure; George Biros, Alex Gray, Rich Vuduc*

## **Lawrence Berkeley National Laboratory**

– *Sam Williams, Lenny Oliker*

Los Alamos Computer Science Symposium 2009

# Key Ideas and Findings

- ▶ First cross-platform single-node multicore study of tuning the fast multipole method (FMM)
  - ▶ Explores data structures, SIMD, multithreading, mixed-precision, and tuning
  - ▶ Show 25x speedups on Intel Nehalem, 9.4x AMD Barcelona, 37.6x Sun Victoria Falls
  - ▶ Surprise? Multicore ~ GPU in performance & energy efficiency for the FMM
- ▶ Broader context: Generalized n-body problems, for particle simulation & statistical data analytics

- ▶ *Context:*  
Mathematical & programming model foundations  
for generalized n-body problems (GNP)

# Context: Interaction calculations

- ▶ What do these have in common?

$$\forall q \in Q : \quad F(q) = \sum_{r \in (Q - \{q\})} C \frac{r - q}{\|r - q\|^3}$$

*Force computation*

$$\forall q \in Q : \quad \text{AllNN}(q) = \operatorname{argmin}_{r \in R} d(q, r)$$

*All nearest neighbors*

$$\forall q \in Q : \quad \text{KDE}(q) = \frac{1}{|R|} \sum_{r \in R} K(q, r)$$

*Kernel density estimation*

$$\forall q \in Q : \quad \text{Range}(q) = \sum_{r \in R} I(\text{dist}(q, r)) \leq h$$

*Range count*

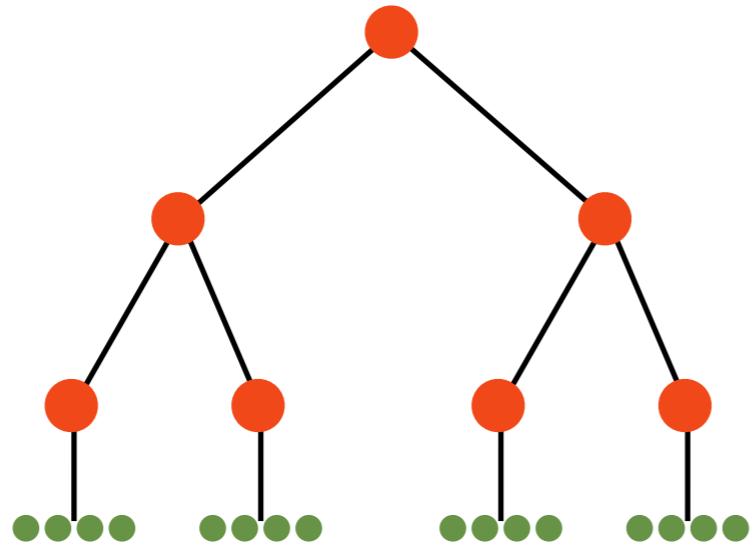
- ▶ Consider pairs of points – naively  $\mathcal{O}(N^2)$

# Commonality: Optimal approximation algorithms

$$\forall q \in Q : \quad F(q) = \sum_{r \in (Q - \{q\})} C \frac{r - q}{\|r - q\|^3}$$

*Force computation*

- ▶ Hierarchical tree-based approximation algorithms for force computations, e.g., Barnes-Hut or FMM



Evaluate interactions  
→ Tree traversals

Store aggregate data at  
nodes, e.g., bounding box,  
mass

# Programming model?

$$\forall q \in Q : \quad F(q) = \sum_{r \in (Q - \{q\})} C \frac{r - q}{\|r - q\|^3}$$

*Force computation*

# Programming model?

$$\forall q \in Q : \quad F(q) = \sum_{r \in (Q - \{q\})} C \frac{r - q}{\|r - q\|^3}$$

*Force computation*

- ▶ Write in a more suggestible form

$$\Phi(Q) \equiv \text{map}_{q \in Q} \sum_{r \in (Q - \{q\})} C \frac{r - q}{\|r - q\|^3}$$

*MapReduce-like*

# Programming model: THOR

$$\forall q \in Q : \quad F(q) = \sum_{r \in (Q - \{q\})} C \frac{r - q}{\|r - q\|^3}$$

*Force computation*

- ▶ Write in a more suggestible form

$$\Phi(Q) \equiv \text{map}_{q \in Q} \sum_{r \in (Q - \{q\})} C \frac{r - q}{\|r - q\|^3}$$

*MapReduce-like*

- ▶ Idea: Tree-based High-Order Reduce (THOR)
  - ▶ Provide MapReduce-like programming model
  - ▶ Implementations use **fast** algorithms

# Generalized n-body problems (GNPs)

- ▶ General “query-reference” form:

$\forall q_1 \in Q_1, \dots, q_I \in Q_I :$

$$\text{GNP}(q_1, \dots, q_I) = g(q_1, \dots, q_I, \bigotimes_{r_1 \in R_1} \dots \bigotimes_{r_J \in R_J} f(q_1, \dots, q_I, r_1, \dots, r_J))$$

- ▶ Simple form,  $I = J = I$ :

$$\forall q \in Q : \text{GNP}(q) = g(q, \bigotimes_{r \in R} f(q, r))$$

- ▶ End-user programmer fills in **green** stuff

# Main research questions

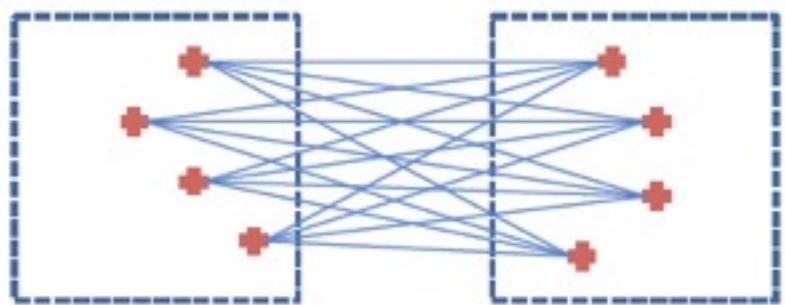
- ▶ **[Math/algorithms]** Optimal approximation algorithms with accuracy guarantees?
- ▶ **[Programming models]** How to best express such methods at a “high-level?”
- ▶ **[Implementation]** How to translate from high-level to efficient low-level?
  - ▶ Big data motivates fast (optimal) algorithms
  - ▶  $O(N)$  is good, but constant matters!
  - ▶ Best architectures? Data structures? Numerics? Tuning?

- ▶ High-performance multicore FMMs:  
Analysis, optimization, and tuning
  - ▶ Algorithmic characteristics
  - ▶ Architectural implications
  - ▶ Observations

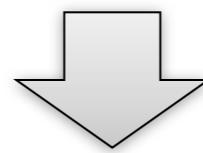
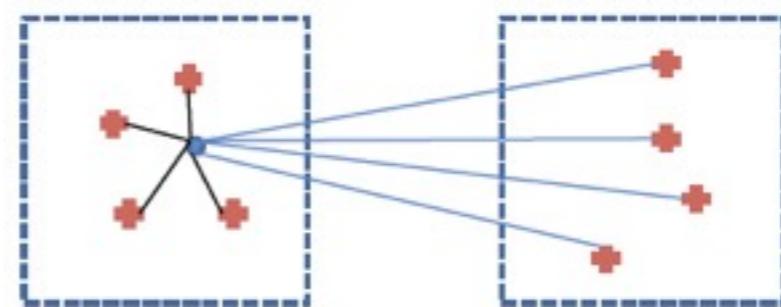
- ▶ High-performance multicore FMMs:  
Analysis, optimization, and tuning
  - ▶ **Algorithmic characteristics**
  - ▶ Architectural implications
  - ▶ Observations

# Computing Direct vs. Tree-based Interactions

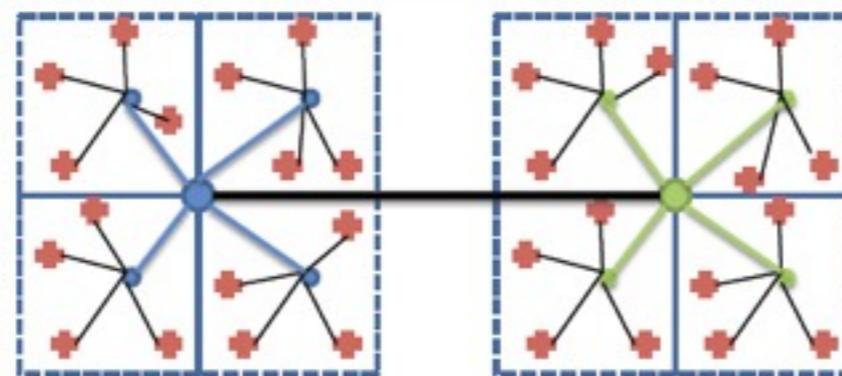
Direct evaluation:  $O(N^2)$



Barnes-Hut:  $O(N \log N)$

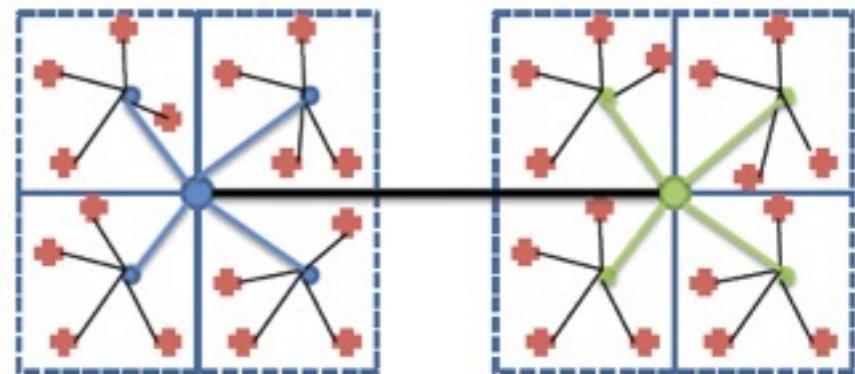


Fast Multipole Method (FMM):  $O(N)$

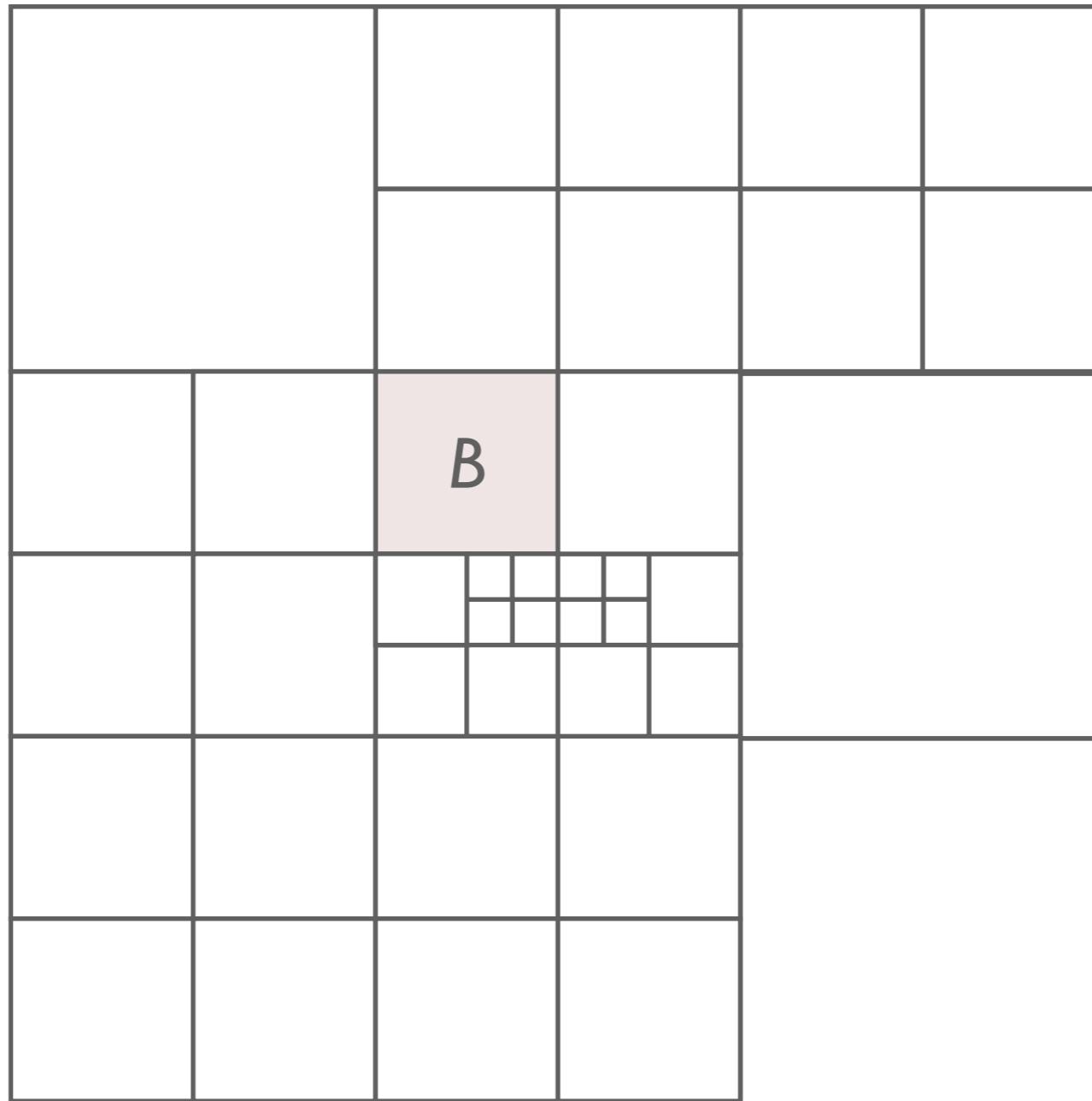


# Fast multipole method

- ▶ Given:
  - ▶  $N$  target points and  $N$  sources
  - ▶ Tree type & max points per leaf,  $q$
  - ▶ Desired accuracy,  $\epsilon$
- ▶ Build tree
- ▶ Evaluate potential at all  $N$  targets
  - ▶ Upward pass, “interact” node with parent
  - ▶ Downward pass, interact parent with children
  - ▶ Perform U-, V-, W-, and X-list interactions

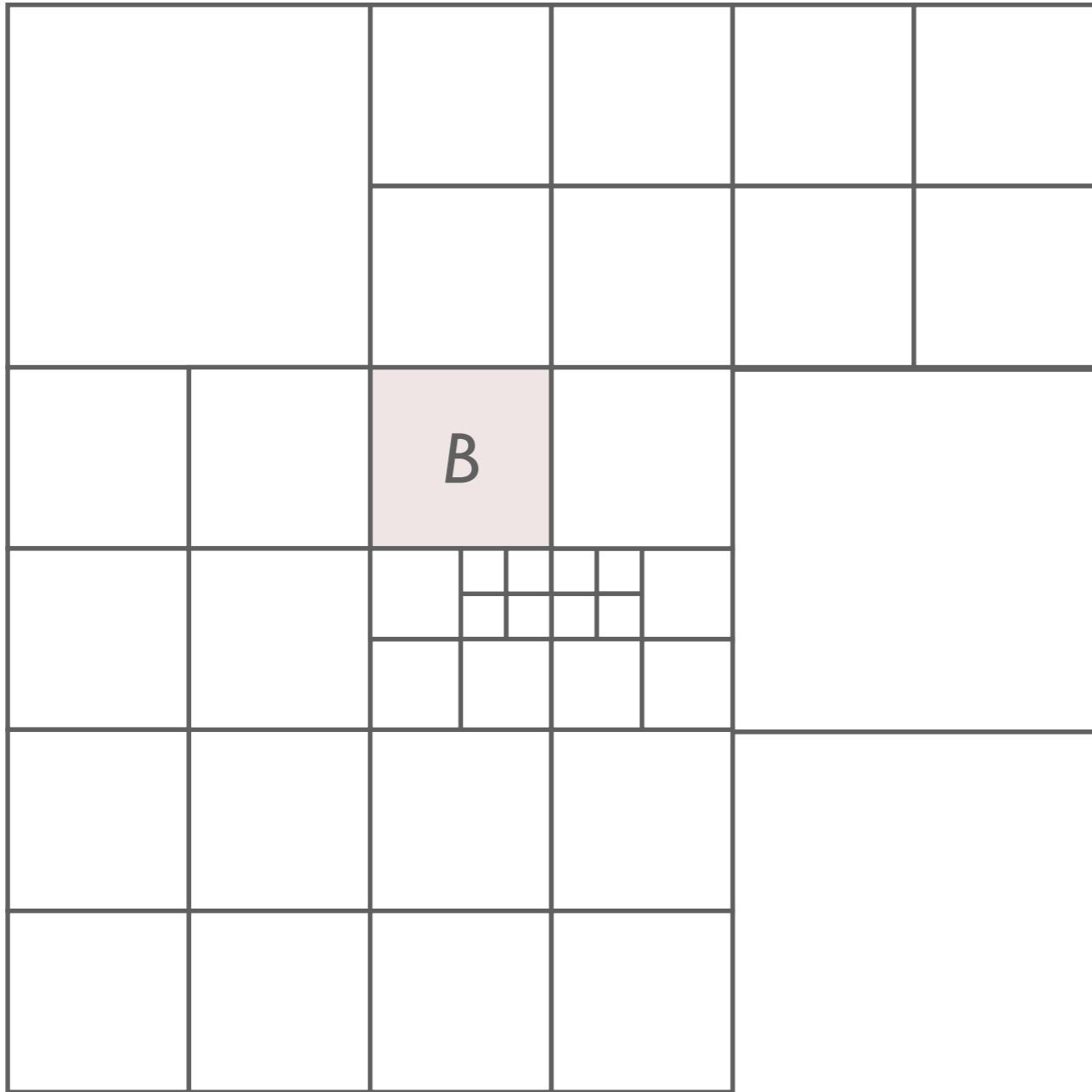


We use *kernel-independent FMM* (KIFMM) of Ying, Zorin, Biros (2004).



## Tree construction

Recursively divide space until each box has **at most  $q$  points**.



## Evaluation phase

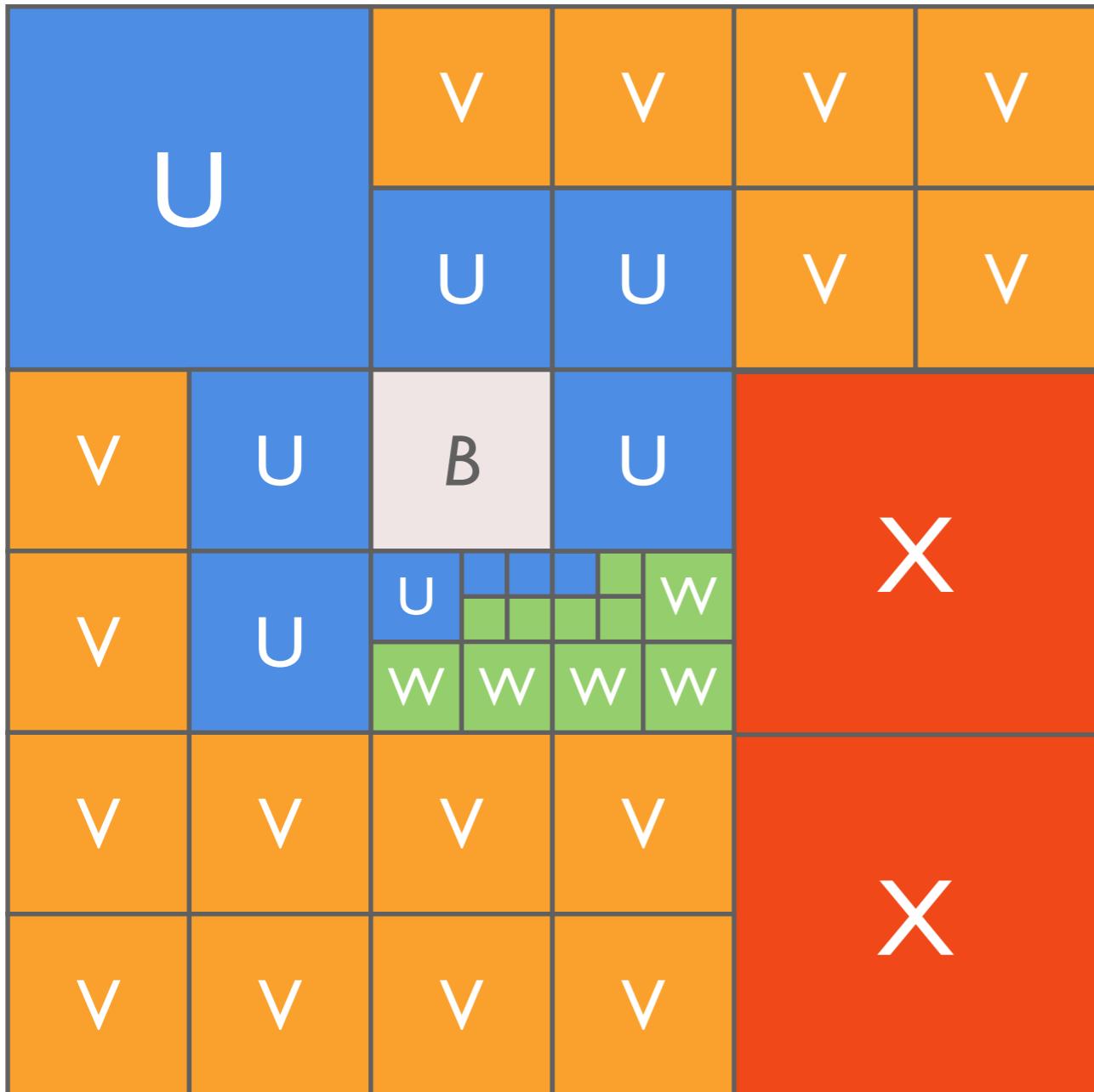
Given the adaptive tree, FMM evaluation performs a series of tree traversals, doing some work at each node,  $B$ .

*Six phases:*

- (1.) Upward pass
- (2–5.) List computations
- (6.) Downward pass

*Phases vary in:*

- data parallelism
- intensity (flops : mops)



## Evaluation phase

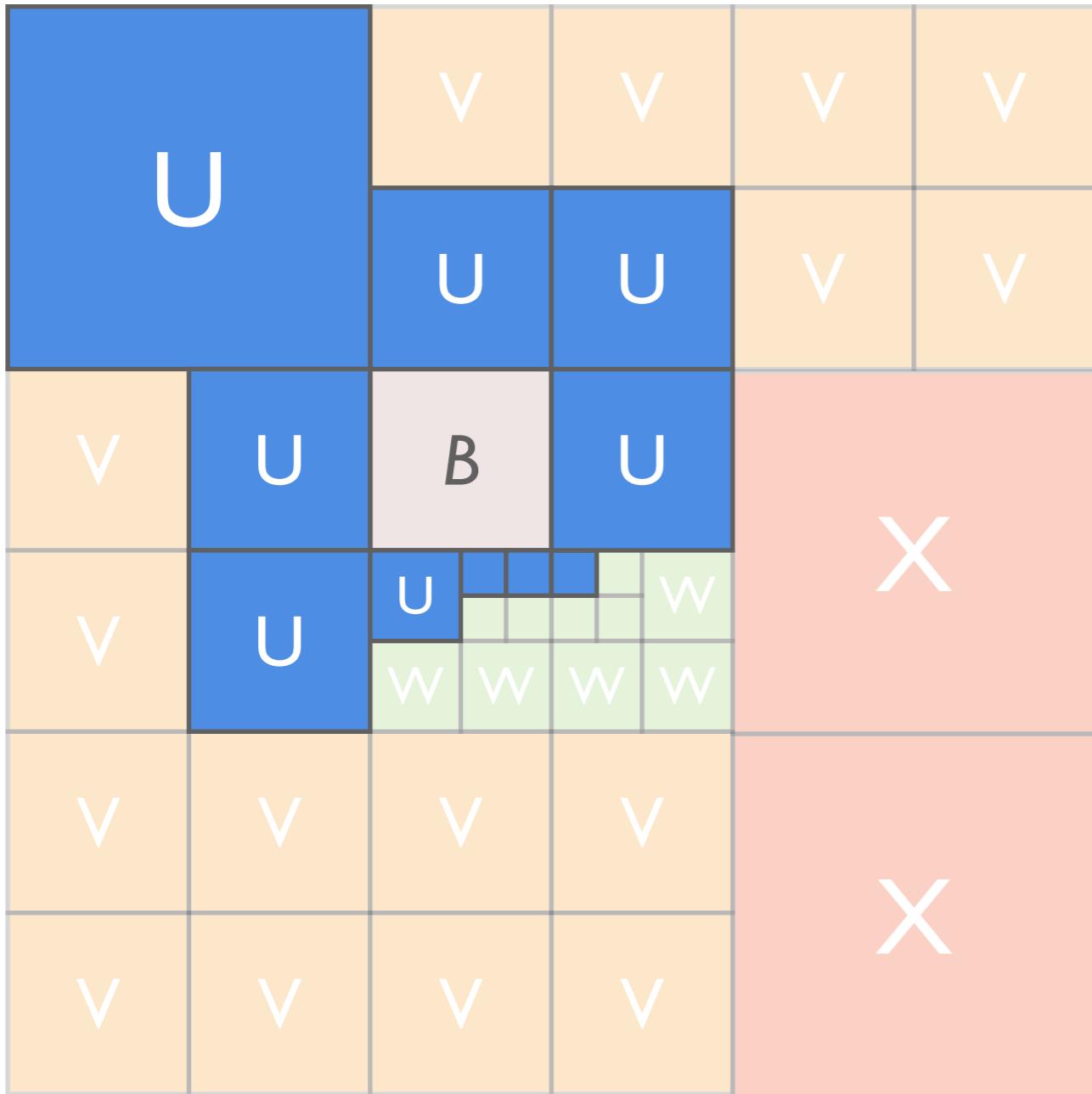
Given the adaptive tree, FMM evaluation performs a series of tree traversals, doing some work at each node,  $B$ .

*Six phases:*

- (1.) Upward pass
- (2–5.) **List computations**
- (6.) Downward pass

*Phases vary in:*

- data parallelism
- intensity (flops : mops)



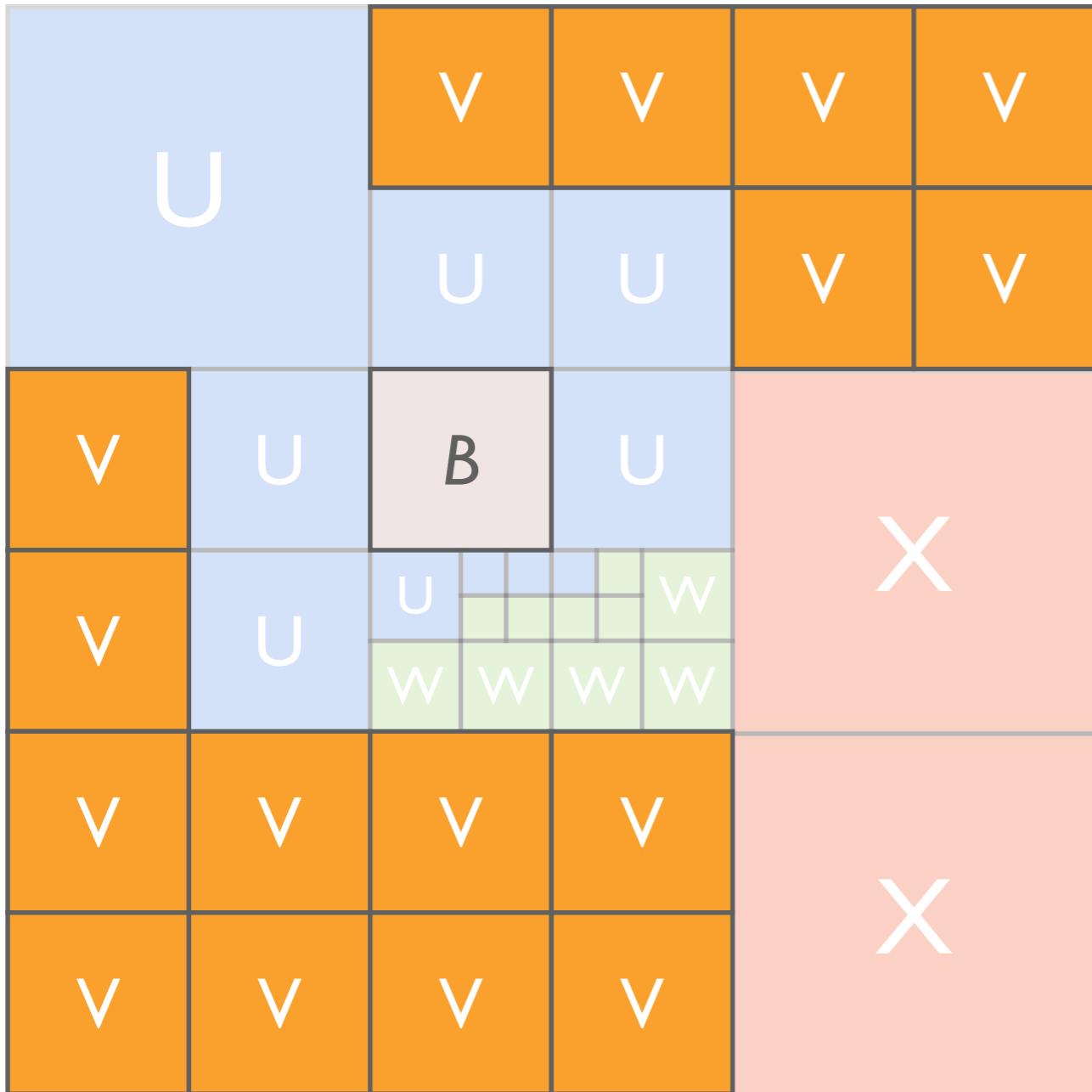
Direct  $B \otimes U$ :

$\rightarrow O(q^2)$  flops :  $O(q)$  mops

## U-List

$U_L(B: leaf) :- \text{neighbors } (B)$

$U_L(B: \text{non-leaf}) :- \text{empty}$

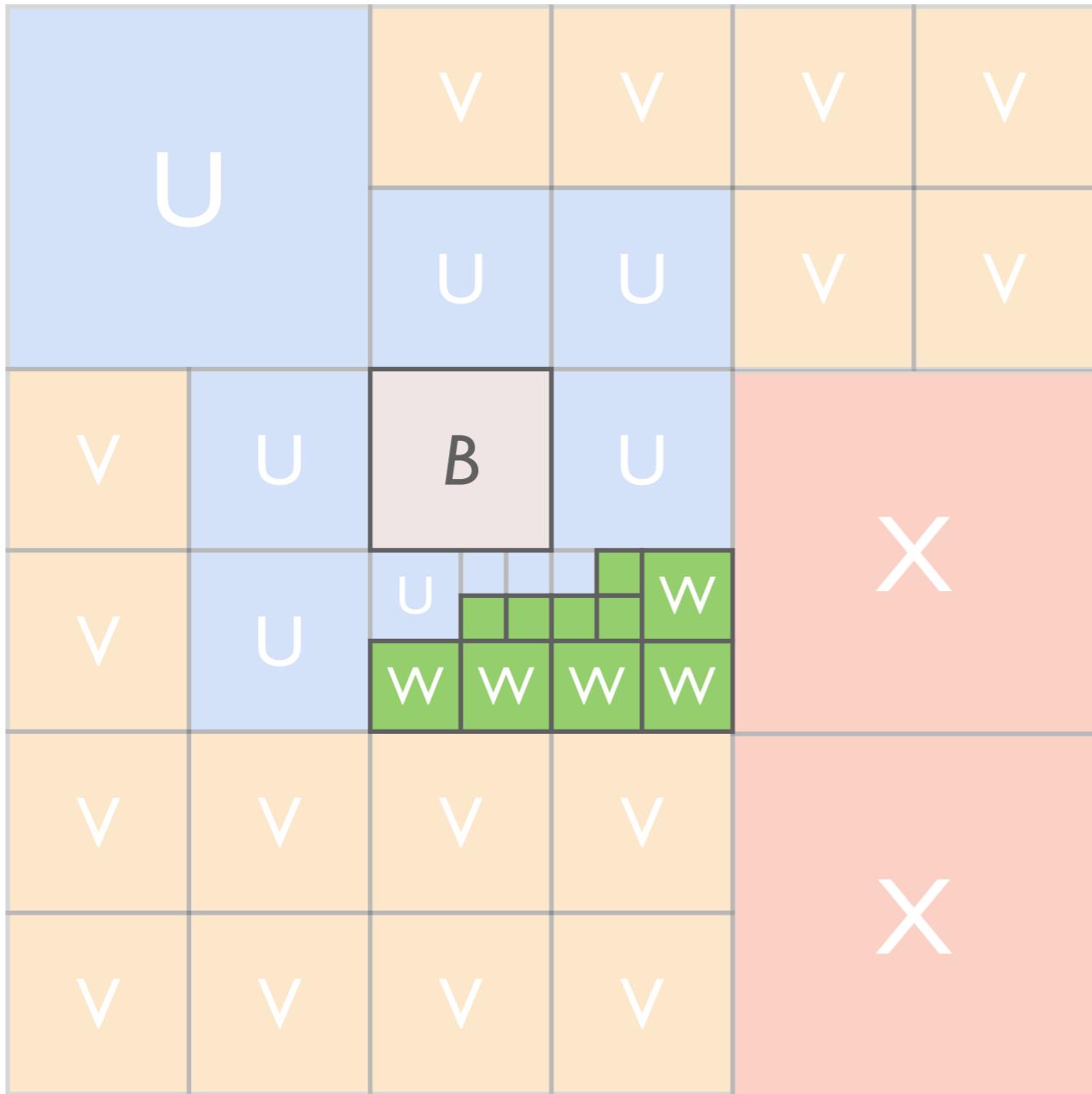


In 3D, FFTs + pointwise multiplication:

- Easily vectorized
- Low intensity vs. U-list

V-List

$V_L(B) := \text{child}(\text{neigh}(\text{par}(B))) - \text{adj}(B)$



Moderate intensity

## W-list

$W_L(B: leaf) :- \text{desc}[\text{par}(\text{neigh}(B)) \cap \text{adj}(B)] - \text{adj}(B)$

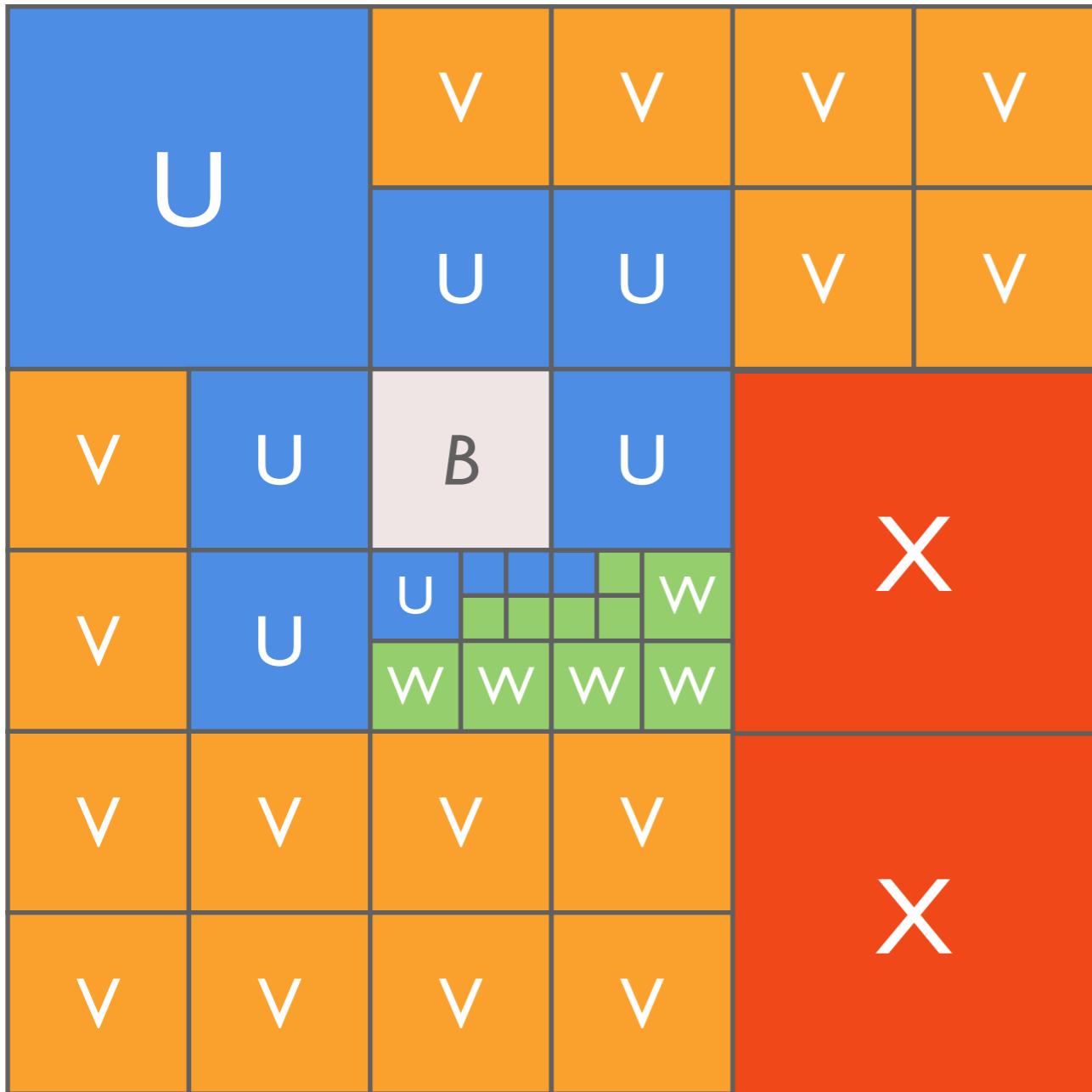
$W_L(B: non-leaf) :- \text{empty}$

		V	V	V
		U	U	V
V	U	B	U	X
V	U	U	W	
V	U	W	W	
V	V	V	V	X
V	V	V	V	X

Moderate intensity

X-list

$X_L(B) :- \{A : B \in W_L(A)\}$

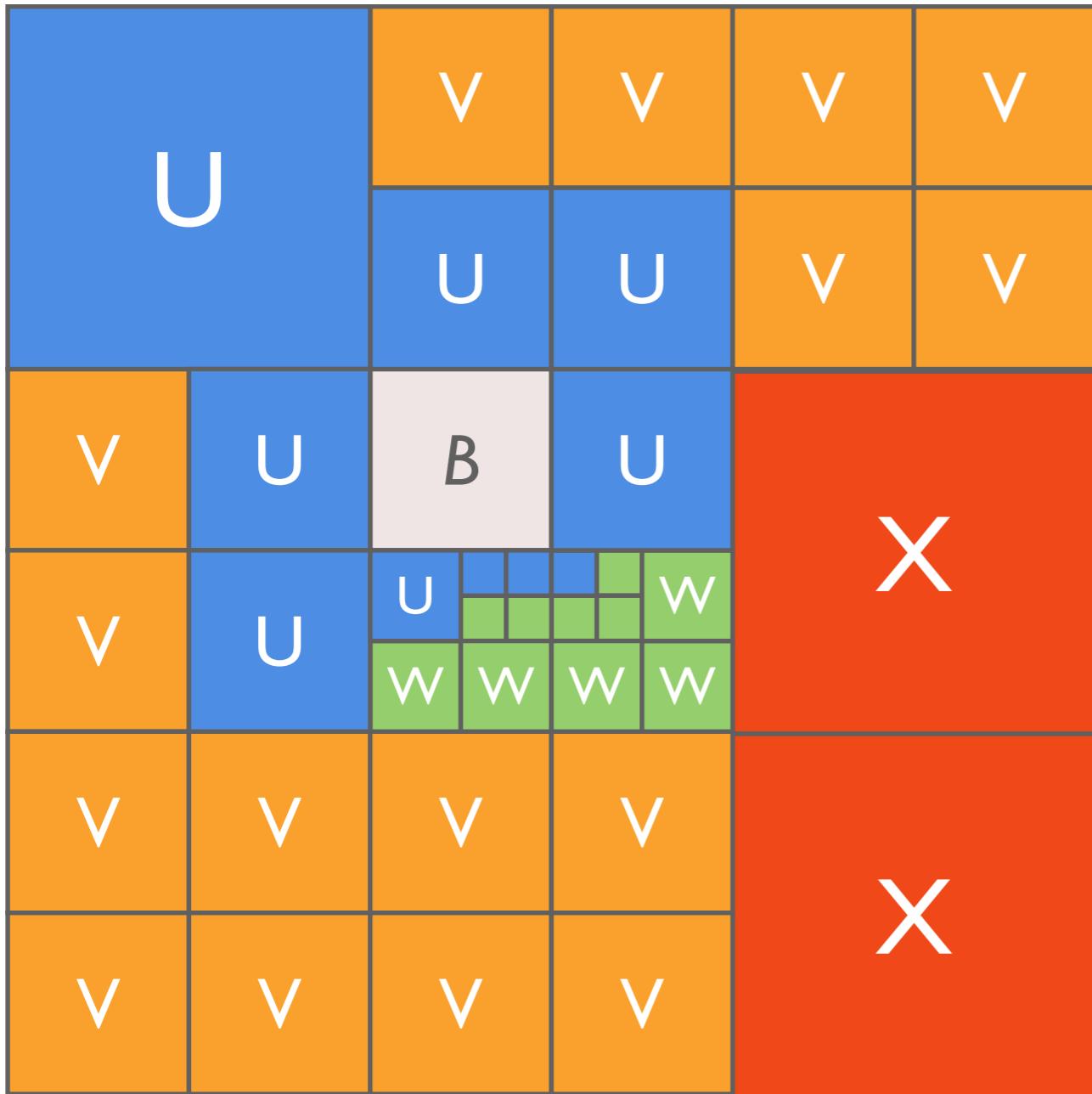


Parallelism exists:

- (1) among phases, with some dependencies;
- (2) within each phase;
- (3) per-box.

*Do not currently exploit (1).*

Essence of the computation



Large  $q$  implies  
 → large U-list cost,  $\mathcal{O}(q^2)$   
 → cheaper V,W,X costs  
 (shallower tree)

Algorithmic tuning  
 parameter,  $q$ , has a global  
 impact on cost.

Essence of the  
 computation

$$K(r) = \frac{C}{\sqrt{r}}$$

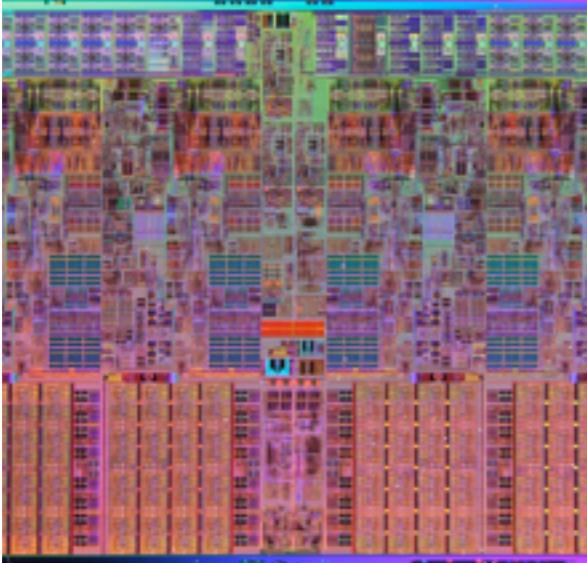
**KIFMM (our variant)**  
requires kernel evaluations  
with expensive flops

## Essence of the computation

*For instance, square-root and divide are expensive, sometimes not pipelined.*

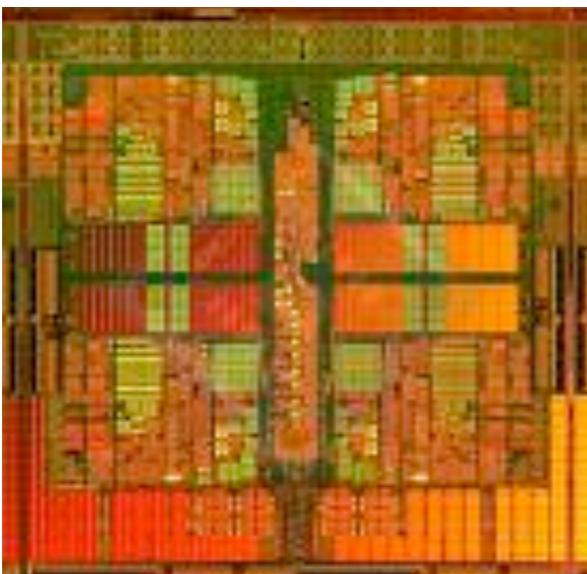
- ▶ High-performance multicore FMMs:  
Analysis, optimization, and tuning
  - ▶ Algorithmic characteristics
  - ▶ **Architectural implications**
  - ▶ Observations

# Hardware thread and core configurations



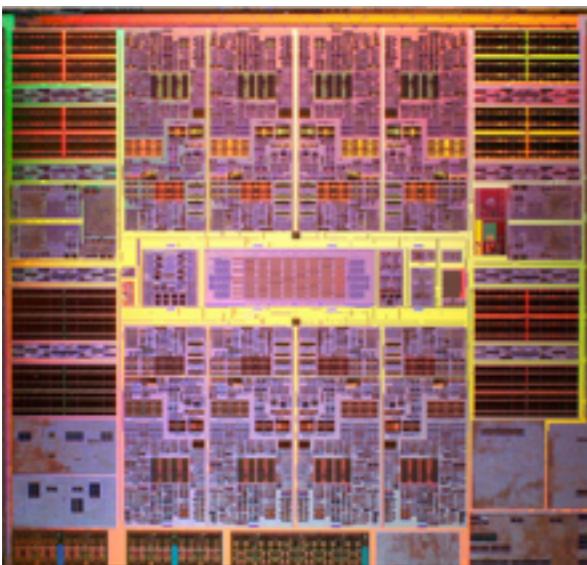
**Intel X5550 “Nehalem”**

(2-socket) × (4-core / socket)



**AMD Opteron 2356 “Barcelona”**

(2-socket) × (4-core / socket)

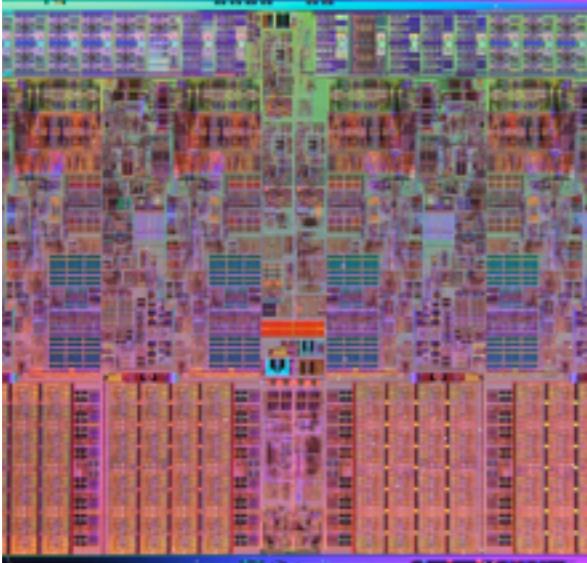


**Sun T5140 “Victoria Falls”**

(2-socket) × (8-core / socket)

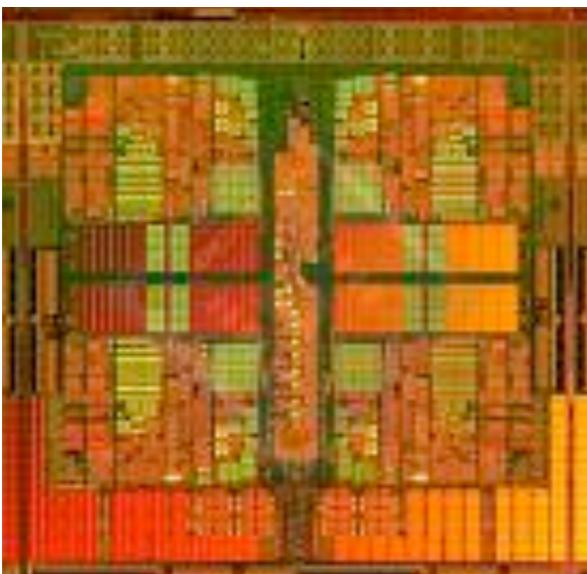
*How do they differ? What implications for FMM?*

# Hardware thread and core configurations



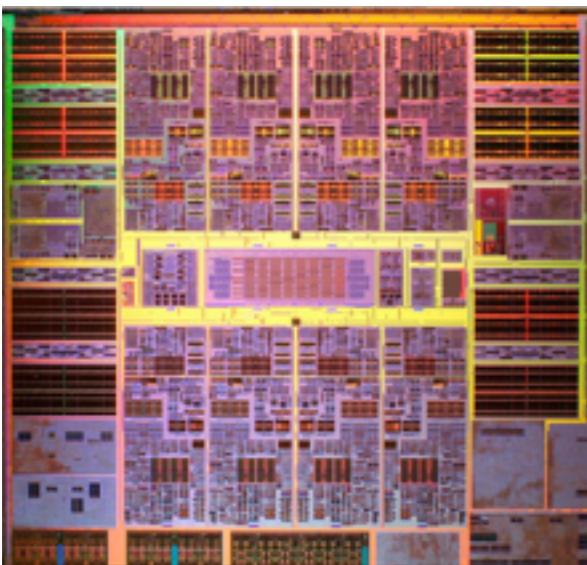
Intel X5550 “Nehalem”

2-sockets  $\times$  4-cores/socket  $\times$  **2-thr/core  $\rightarrow$  16 threads**



AMD Opteron 2356 “Barcelona”

2  $\times$  4  $\times$  **1-thr/core  $\rightarrow$  8 threads**

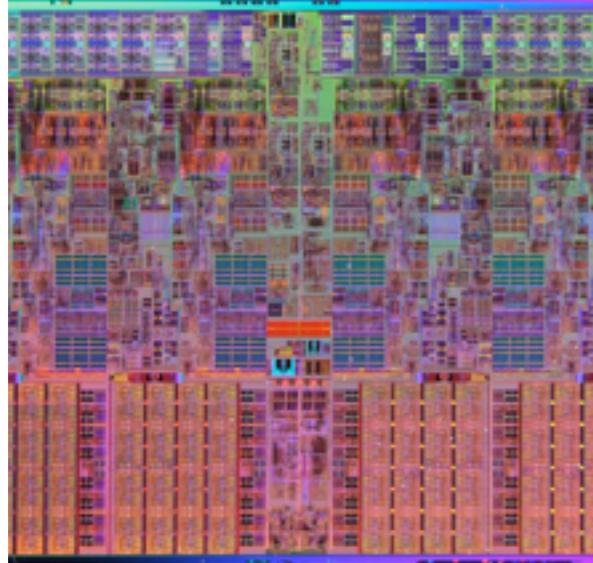


Sun T5140 “Victoria Falls”

2  $\times$  8  $\times$  **8-thr/core  $\rightarrow$  128 threads**

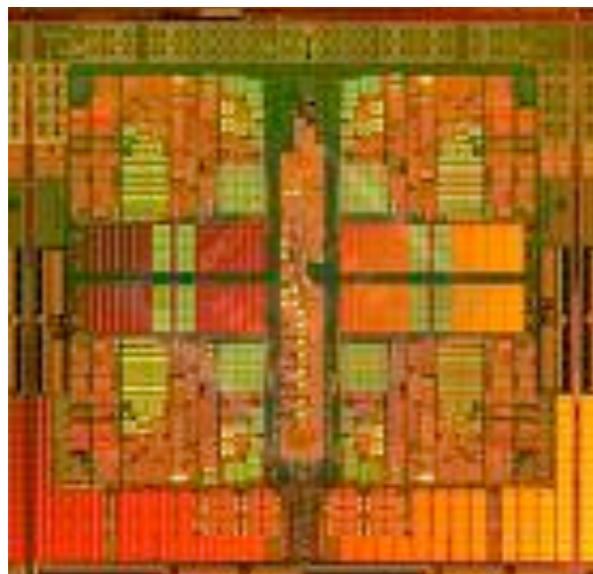
*FMM has fine-grained thread-level parallelism, assuming sufficient functional units.*

## Single core configuration



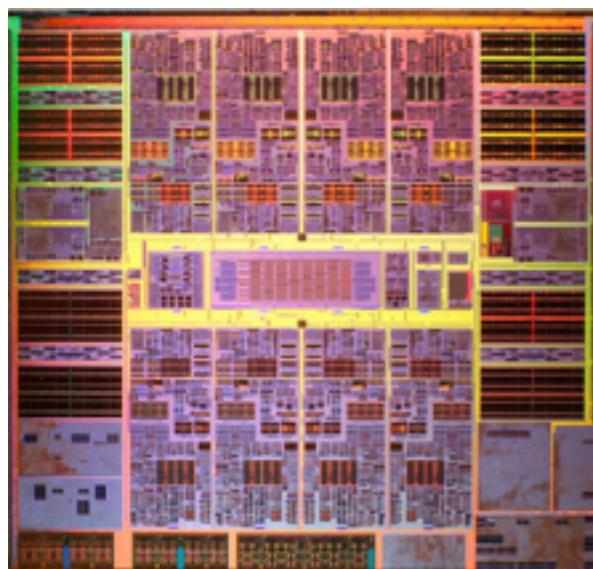
Intel X5550 “Nehalem”

Fast **2.66 GHz** cores, **out-of-order, deep pipelines**.



AMD Opteron 2356 “Barcelona”

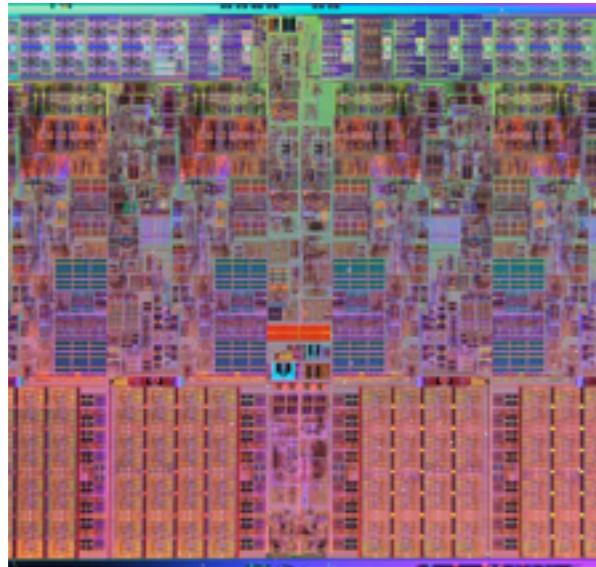
Fast **2.3 GHz** cores, **out-of-order, deep pipelines**.



Sun T5140 “Victoria Falls”

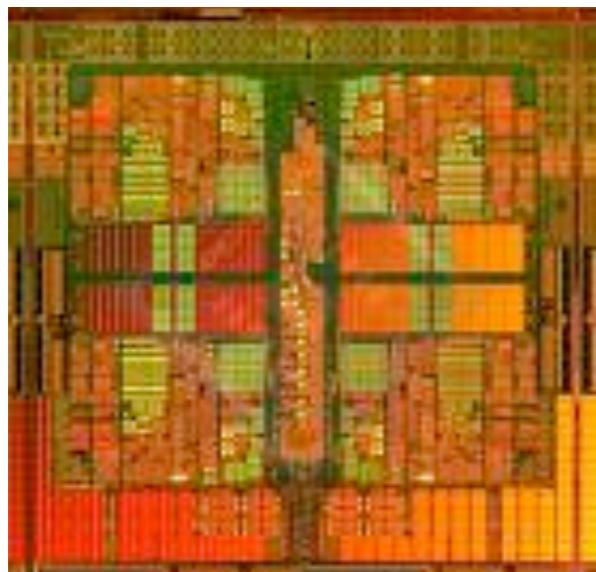
**1.166 GHz** cores, **in-order, shallow pipeline**.

# SIMD



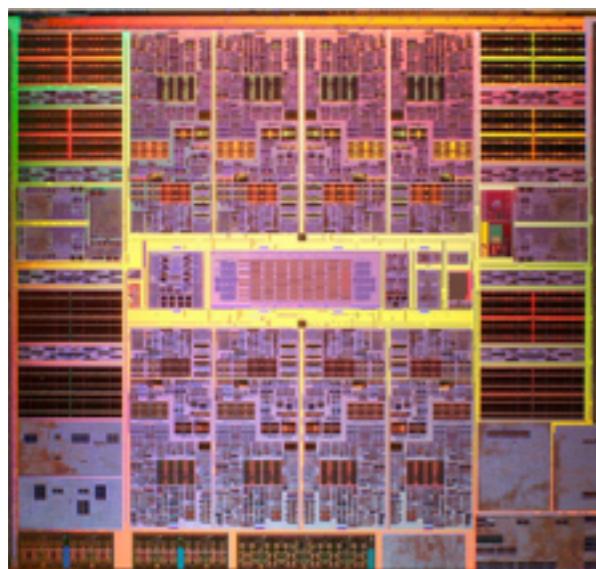
Intel X5550 “Nehalem”

**SIMD** → **85.5** (double), **170.6** (single) **Gflop/s**



AMD Opteron 2356 “Barcelona”

**SIMD** → **73.6** (double), **146.2** (single) **Gflop/s**

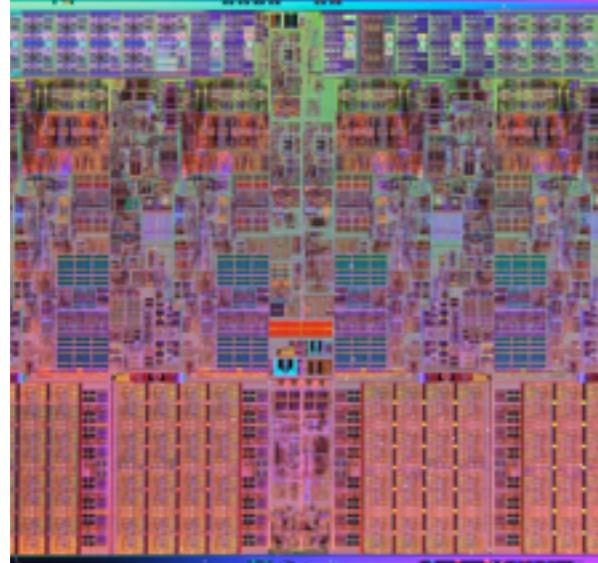


Sun T5140 “Victoria Falls”

**No SIMD** → **18.66 Gflop/s** in single & double

*FMM can use SIMD well, so expect good performance on x86.*

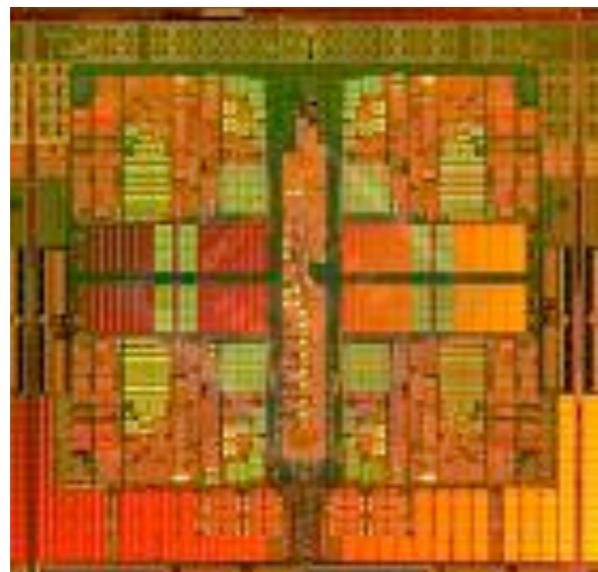
# Floating-point limitations



Intel X5550 “Nehalem”

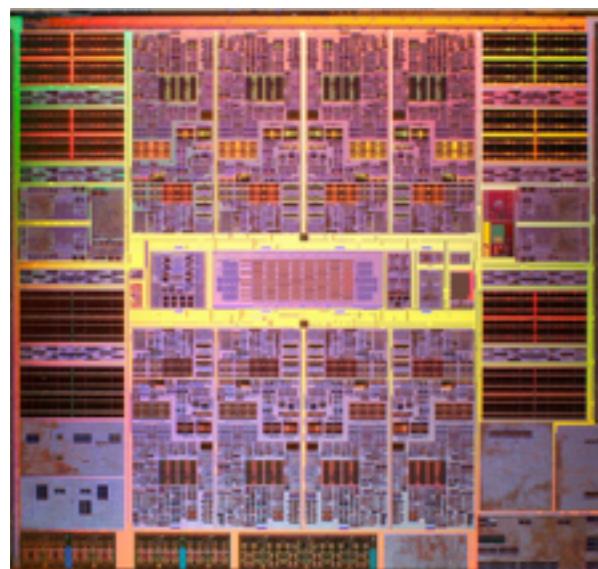
**Reciprocal square-root:**

**0.853** (double), **42.66** (single) **Gflop/s**



AMD Opteron 2356 “Barcelona”

**0.897** (double), **73.6** (single) **Gflop/s**

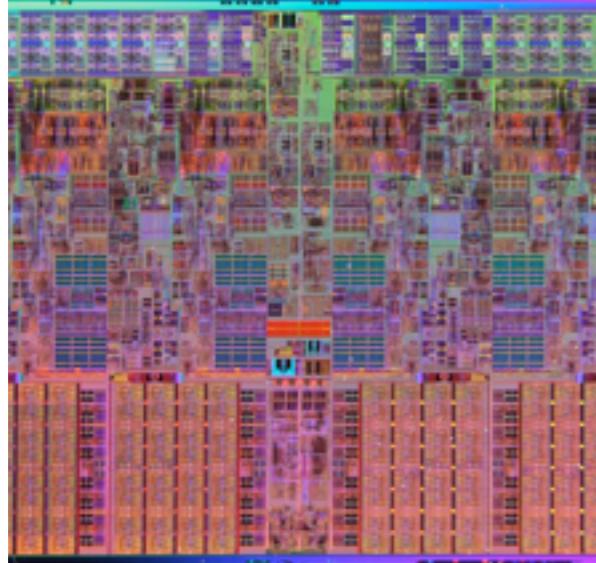


Sun T5140 “Victoria Falls”

**2.26 Gflop/s**

However, x86 has fast approximate single-precision rsqrt, exploitable in double.

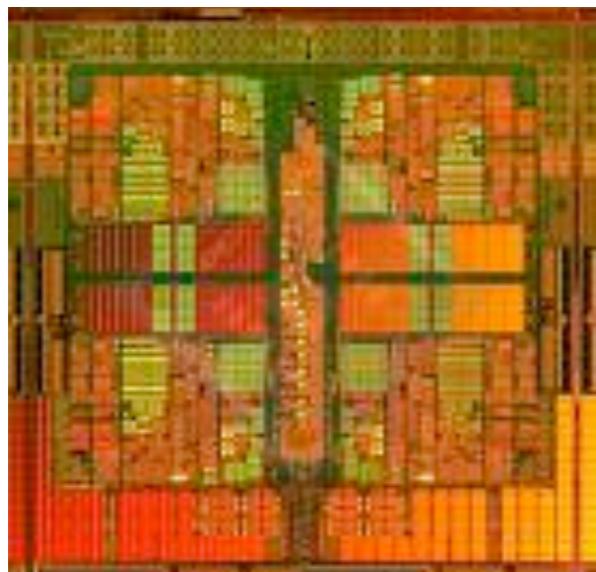
# Memory systems



## Intel X5550 “Nehalem”

Large (**8 MB**) **L3** cache

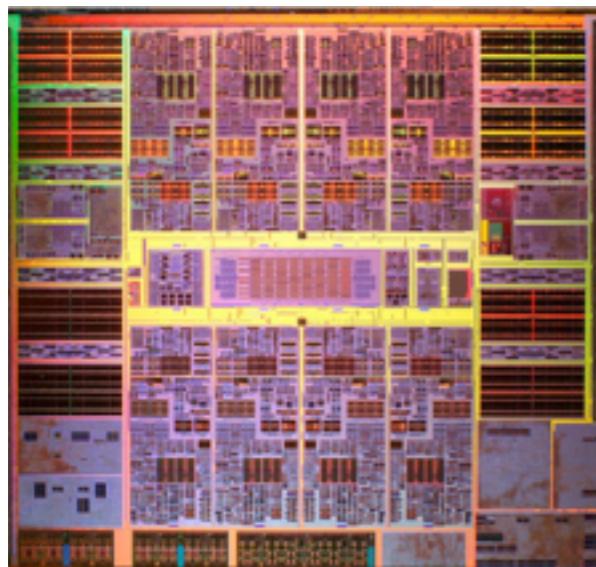
High (**51.2 GB/s**) bandwidth



## AMD Opteron 2356 “Barcelona”

Smaller (**2 MB**) **L3** cache

Lower (**21.33 GB/s**) bandwidth



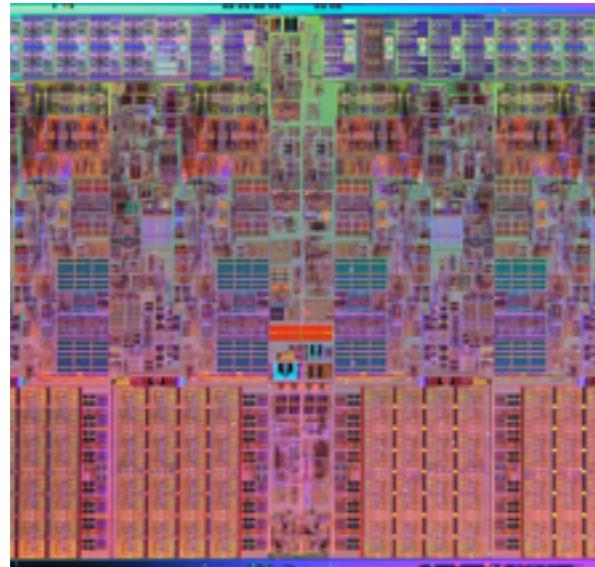
## Sun T5140 “Victoria Falls”

**4 MB L2**

**64.0 GB/s** bandwidth

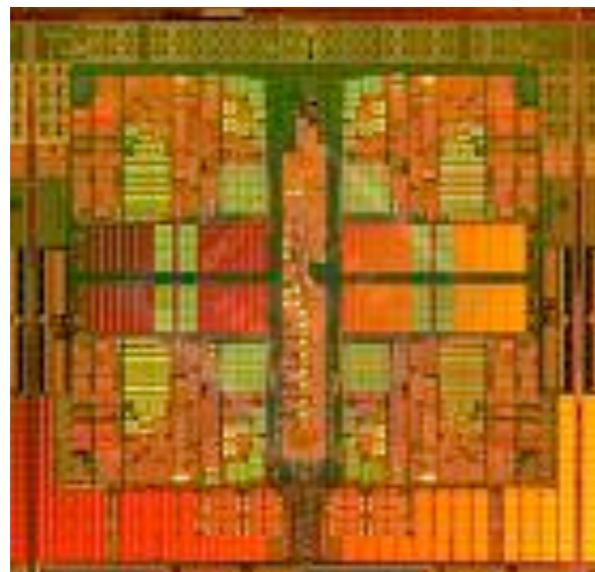
*FMM has a mix of memory behaviors, so memory system impact will vary.*

# Max System Power



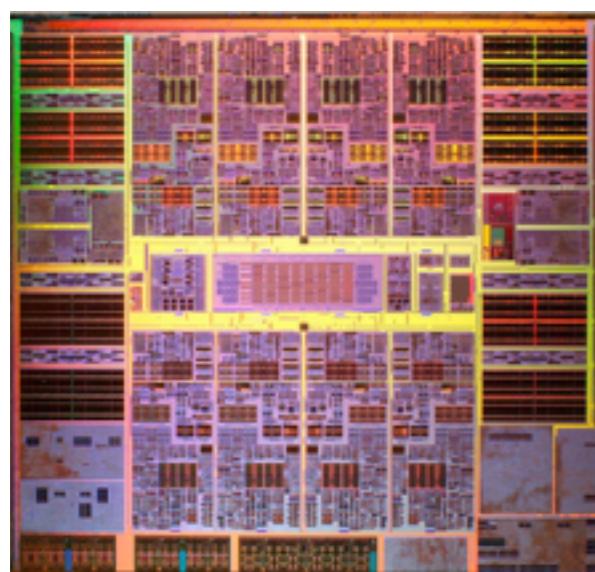
Intel X5550 “Nehalem”

**375 W**



AMD Opteron 2356 “Barcelona”

**350 W**



Sun T5140 “Victoria Falls”

**610 W**

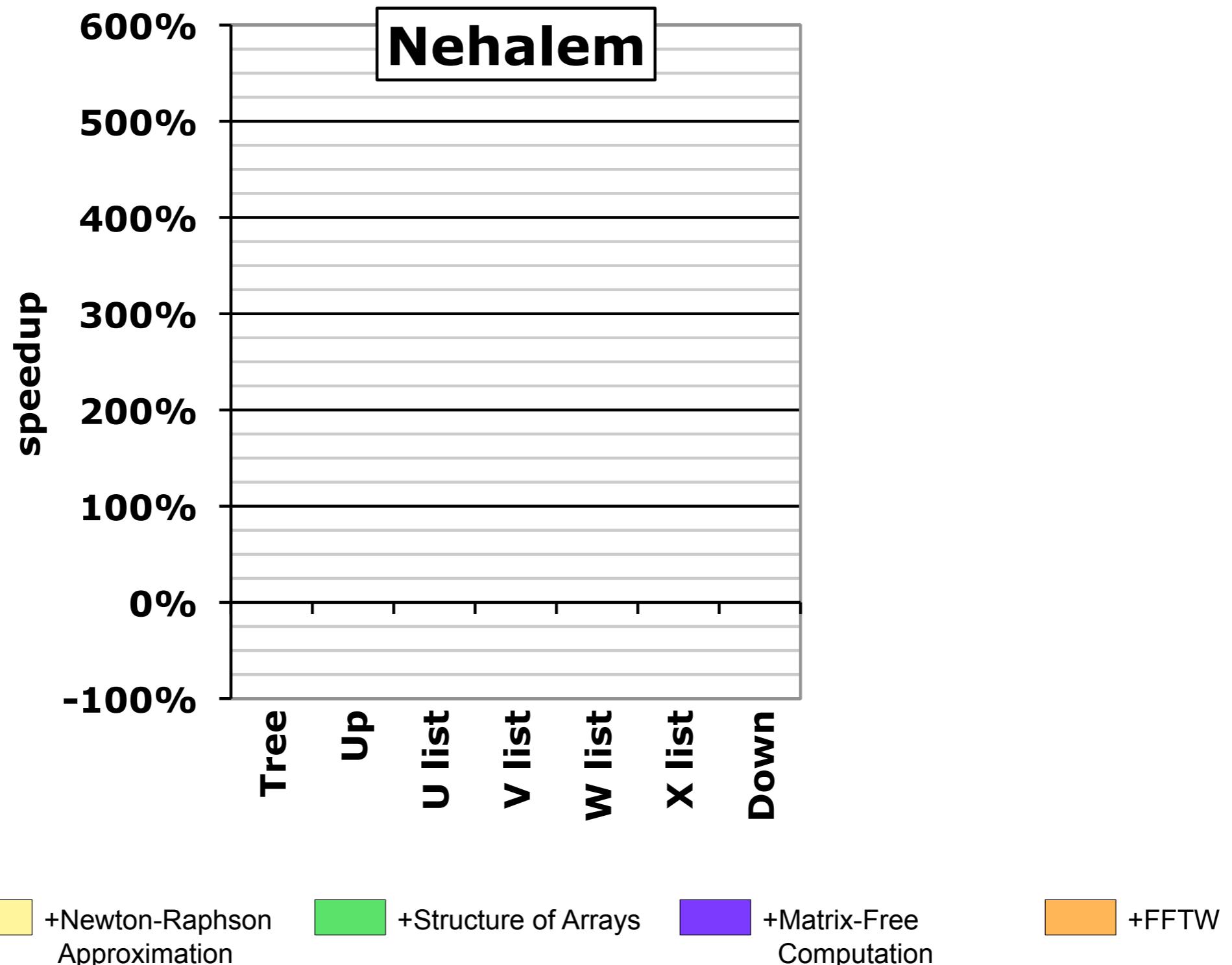
- ▶ High-performance multicore FMMs:  
Analysis, optimization, and tuning
  - ▶ Algorithmic characteristics
  - ▶ Architectural implications
  - ▶ **Observations**

# Optimizations

- ▶ Single-core, manually coded & tuned
  - ▶ *Low-level*: SIMD vectorization (x86)
  - ▶ *Numerical*: rsqrtps + Newton-Raphson (x86)
  - ▶ *Data*: Structure reorg. (transpose or “SOA”)
  - ▶ *Traffic*: Matrix-free via interprocedural loop fusion
  - ▶ FFTW plan optimization
- ▶ OpenMP parallelization
- ▶ Algorithmic tuning of max particles per box,  $q$

# Single-core Optimizations

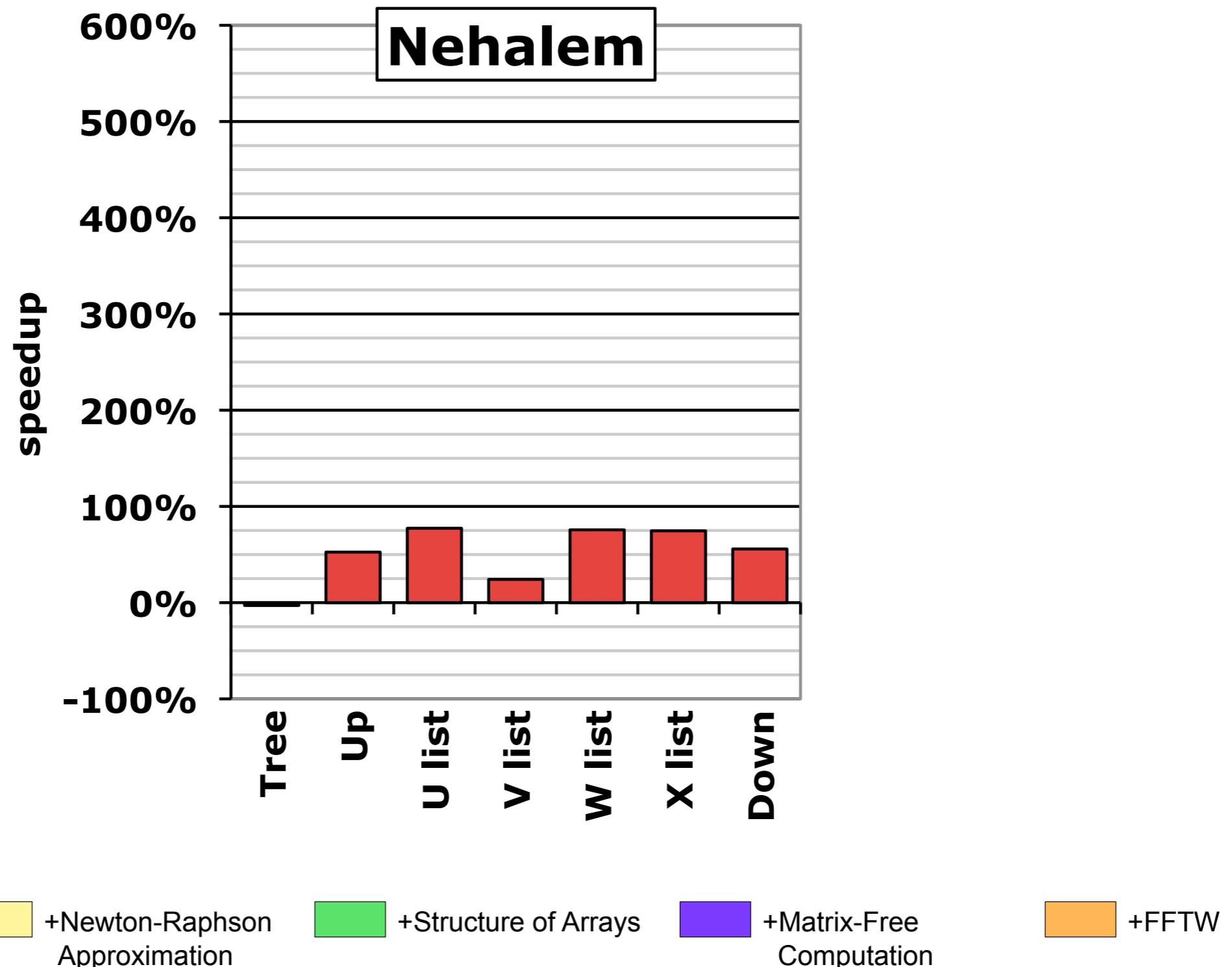
$N_s = N_t = 4M$ , Double-Precision, Non-uniform (ellipsoidal)



Reference: kifmm3d [Ying, Langston, Zorin, Biros]

# Single-core Optimizations

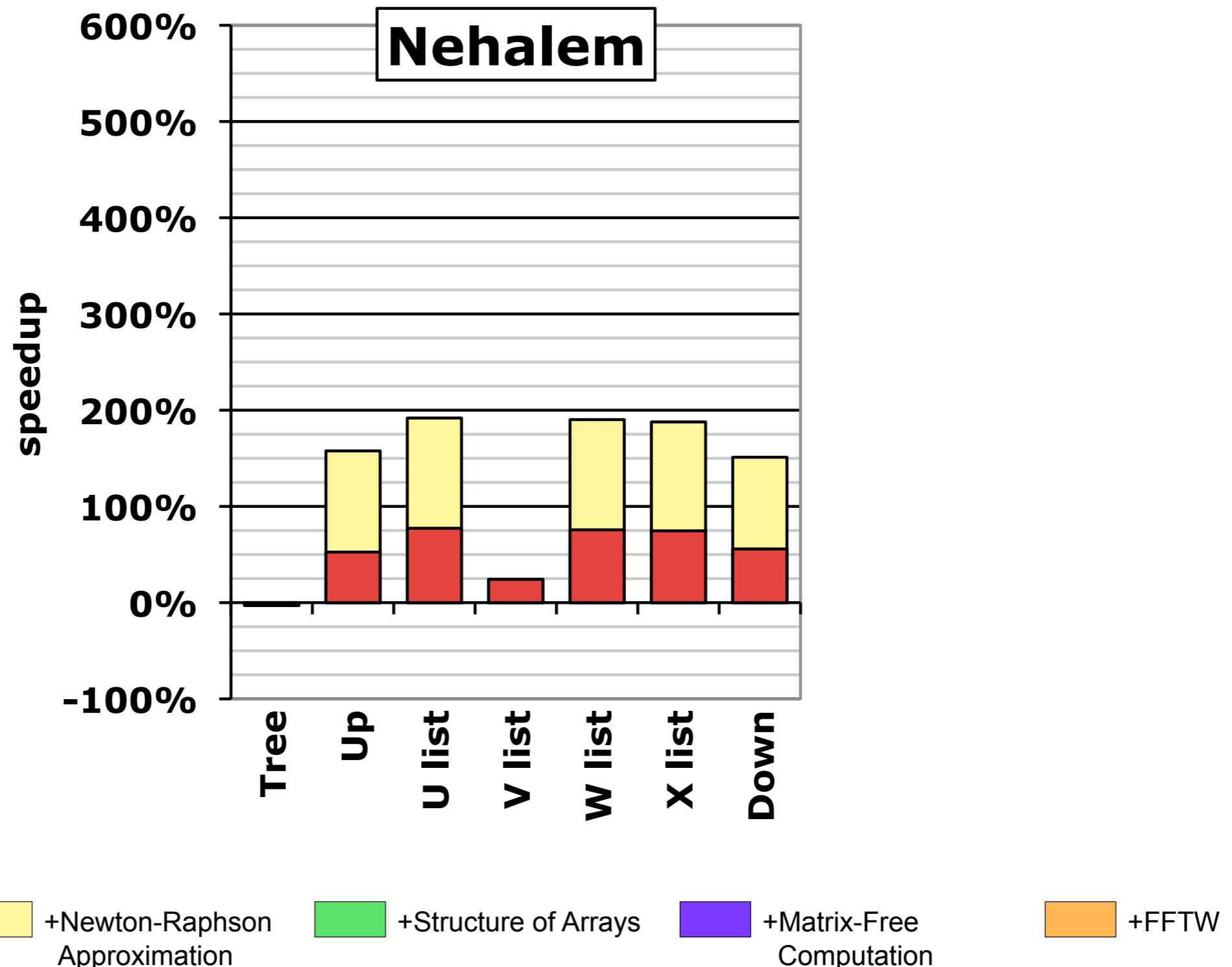
$N_s = N_t = 4M$ , Double-Precision, Non-uniform (ellipsoidal)



SIMD:  $\sim 1.2 - 1.8x$

# Single-core Optimizations

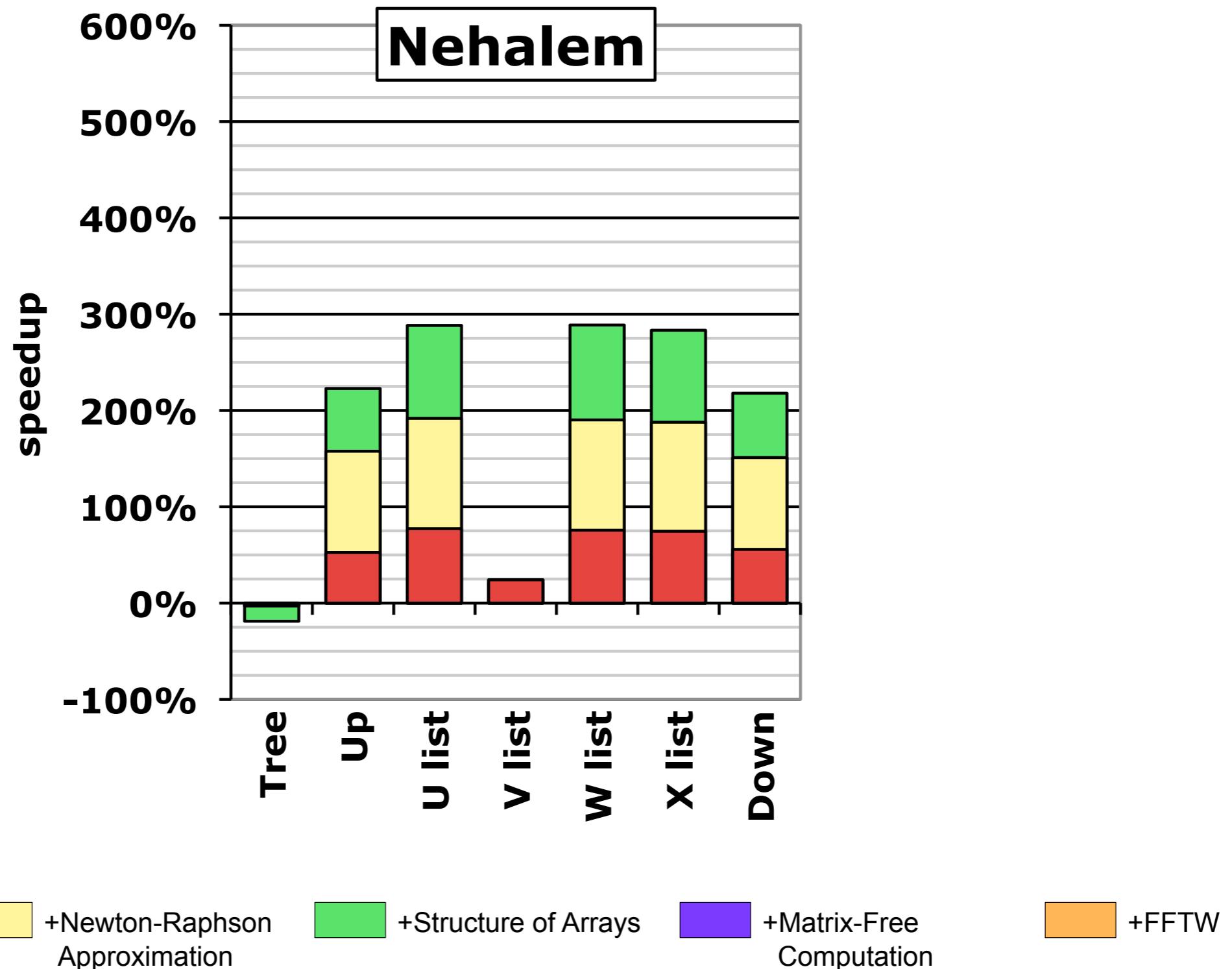
$N_s = N_t = 4M$ , Double-Precision, Non-uniform (ellipsoidal)



*Cumulatively, ~ 3x from SIMD+“smarter numerics”*

# Single-core Optimizations

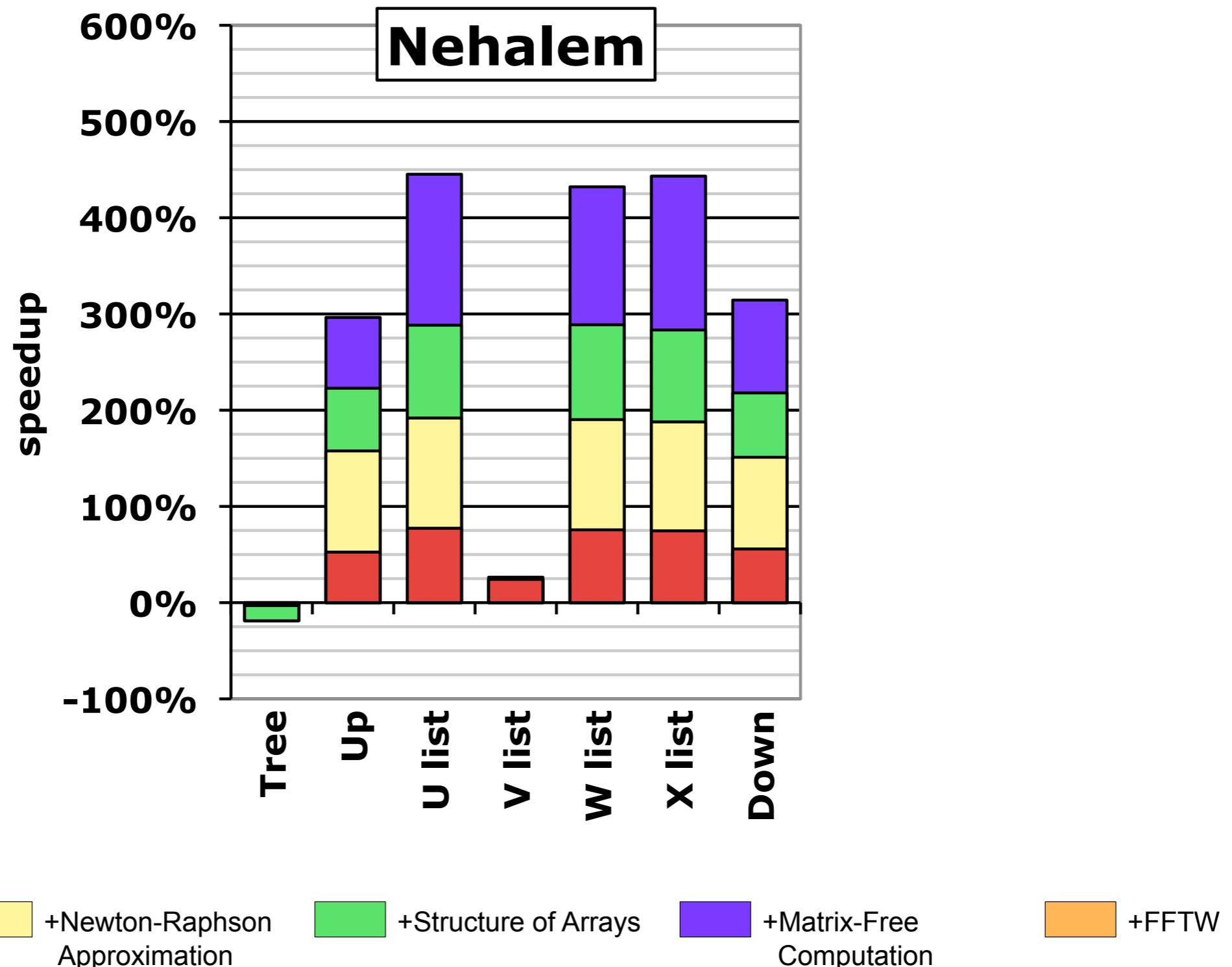
$N_s = N_t = 4M$ , Double-Precision, Non-uniform (ellipsoidal)



Better memory performance + SIMDization, at small cost (“Tree”).

# Single-core Optimizations

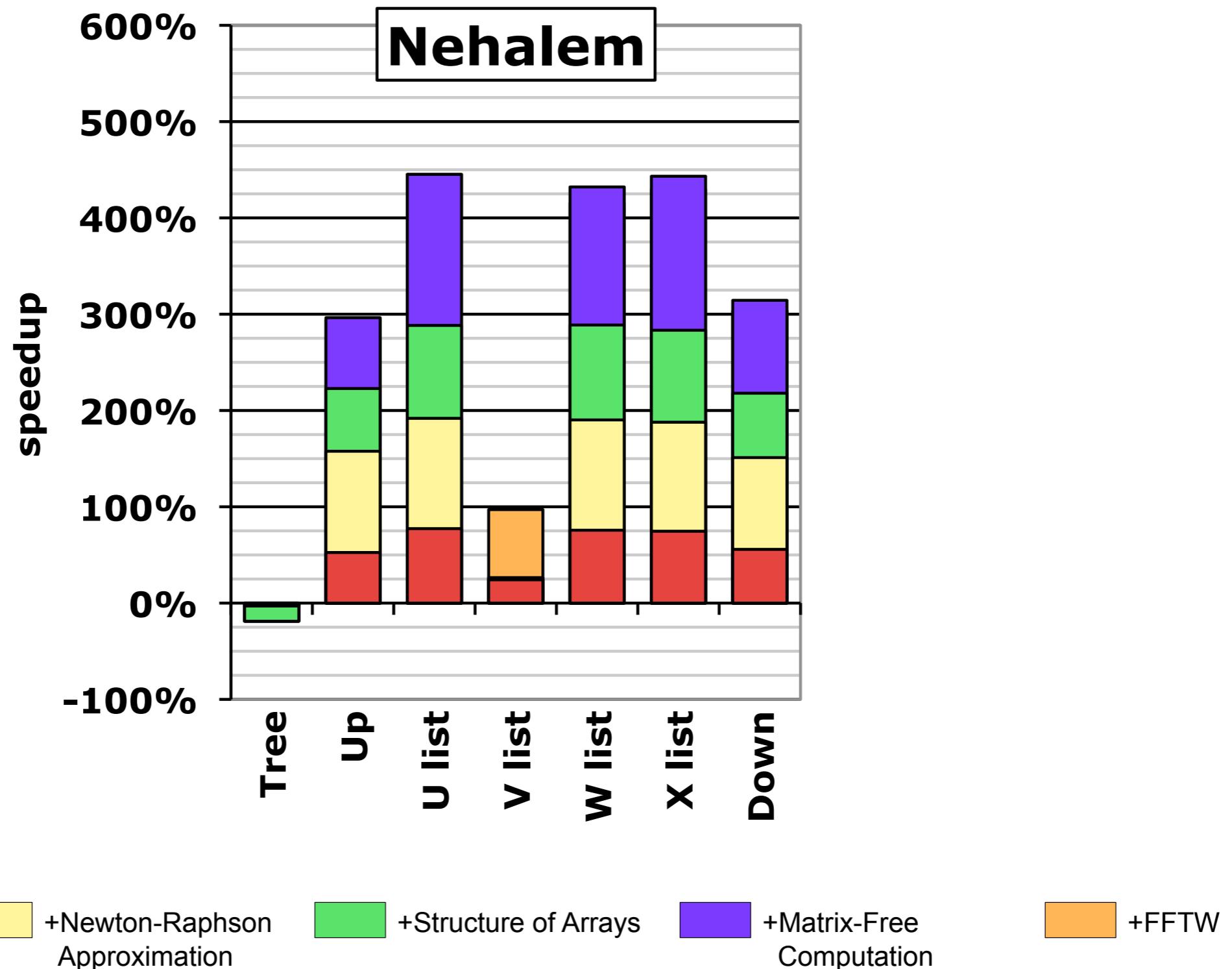
$N_s = N_t = 4M$ , Double-Precision, Non-uniform (ellipsoidal)



SOA+Matrix-free (memory system optimizations)  $\Rightarrow > 2x$

# Single-core Optimizations

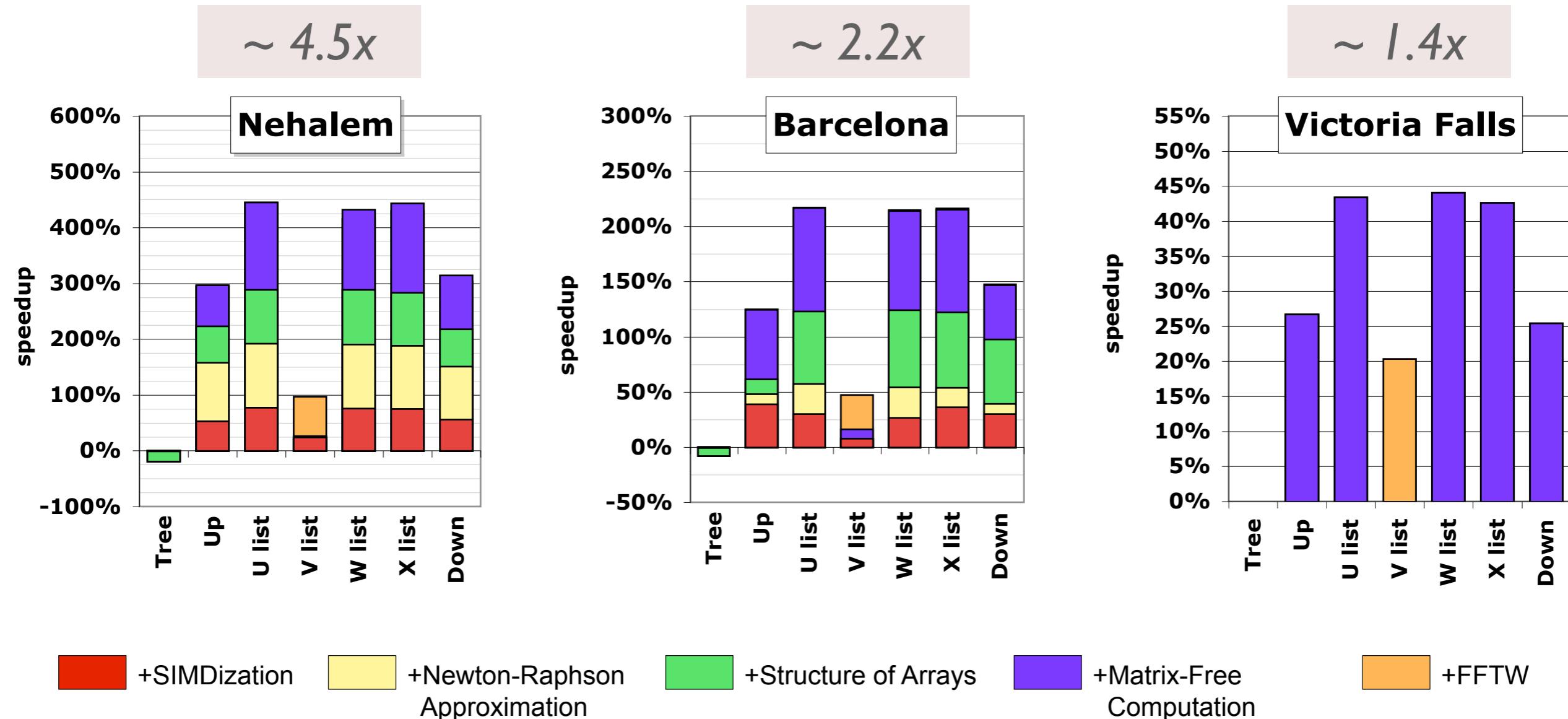
$N_s = N_t = 4M$ , Double-Precision, Non-uniform (ellipsoidal)



“Trick” to improve reduce cost of FFTW plan construction

# Single-core Optimizations

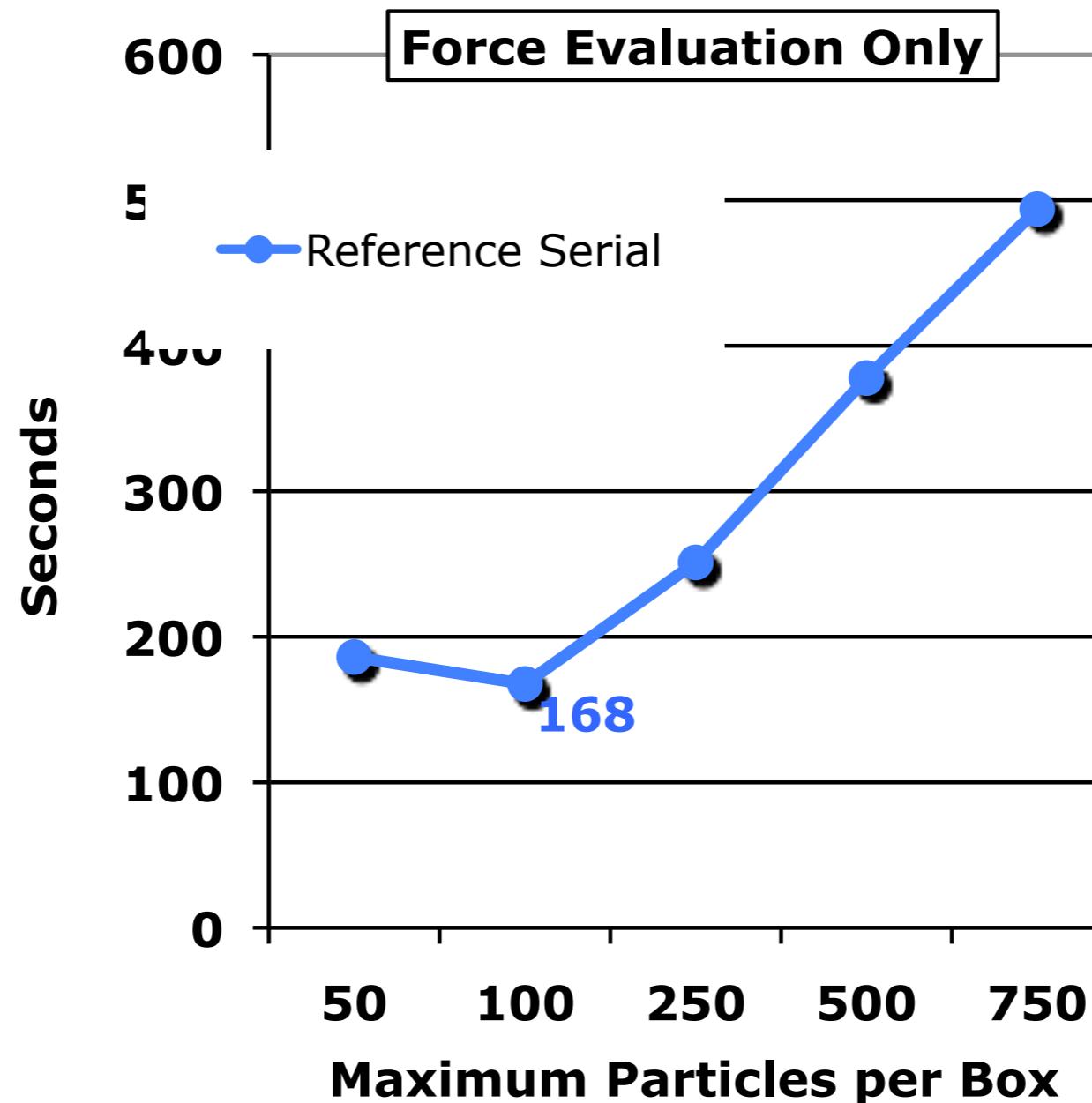
$N_s = N_t = 4M$ , Double-Precision, Non-uniform (ellipsoidal)



*Less impact on Barcelona (why?) and Victoria Falls.*

# Algorithmic Tuning of $q = \text{Max pts / box}$

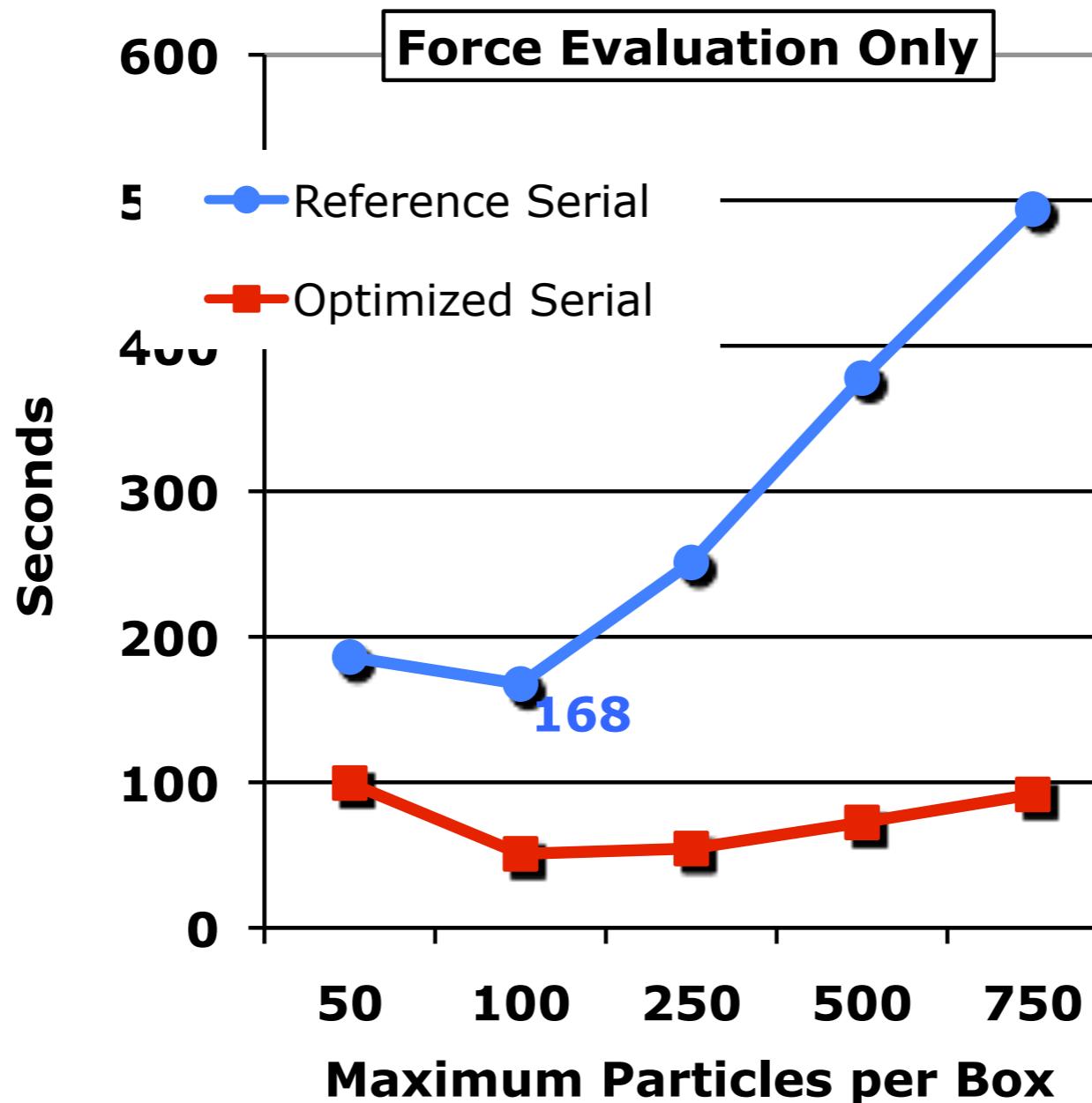
Nehalem



*Tree shape and relative component costs vary as q varies.*

# Algorithmic Tuning of $q = \text{Max pts / box}$

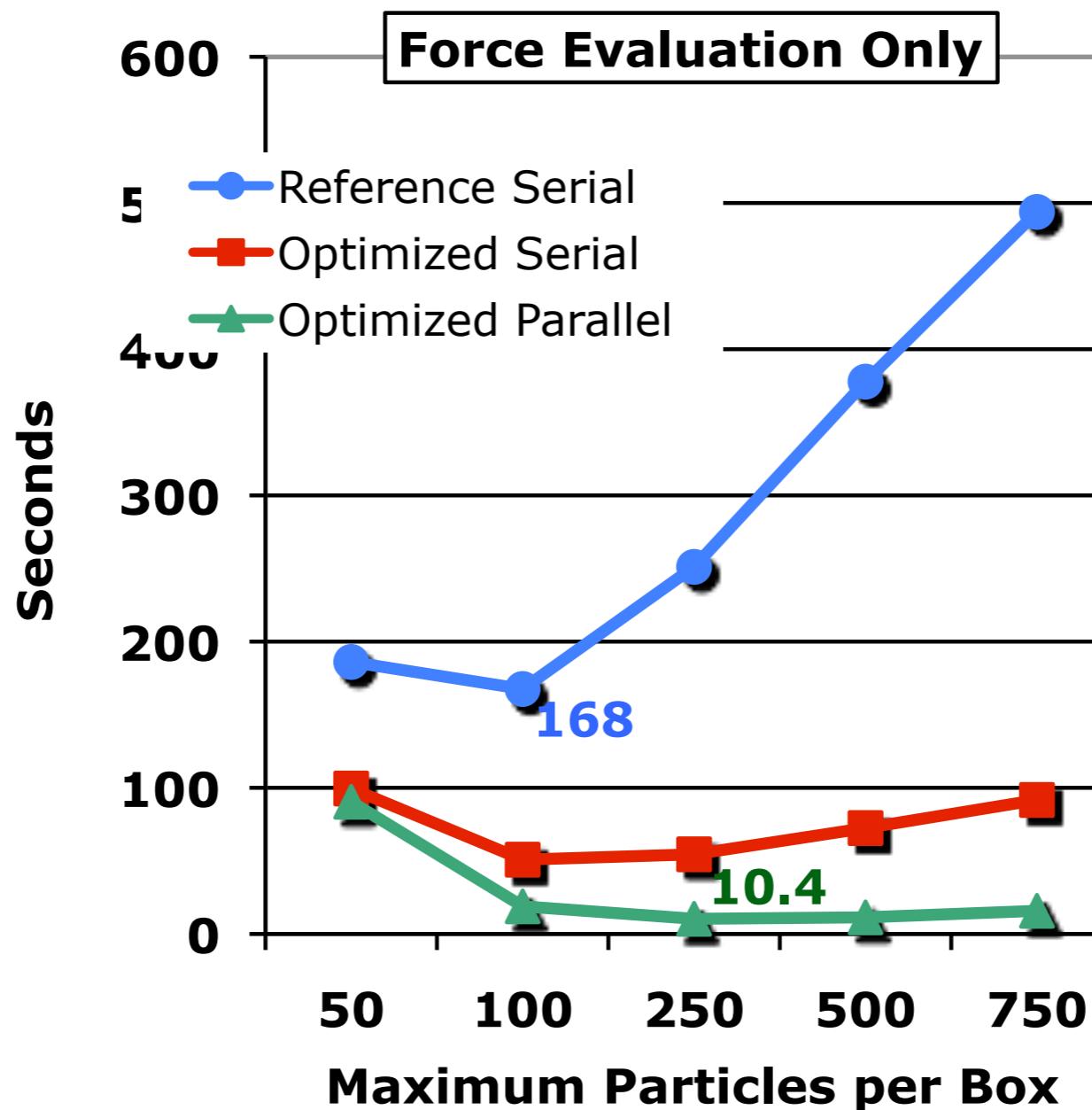
Nehalem



*Shape of curve changes as we introduce optimizations.*

# Algorithmic Tuning of $q = \text{Max pts} / \text{box}$

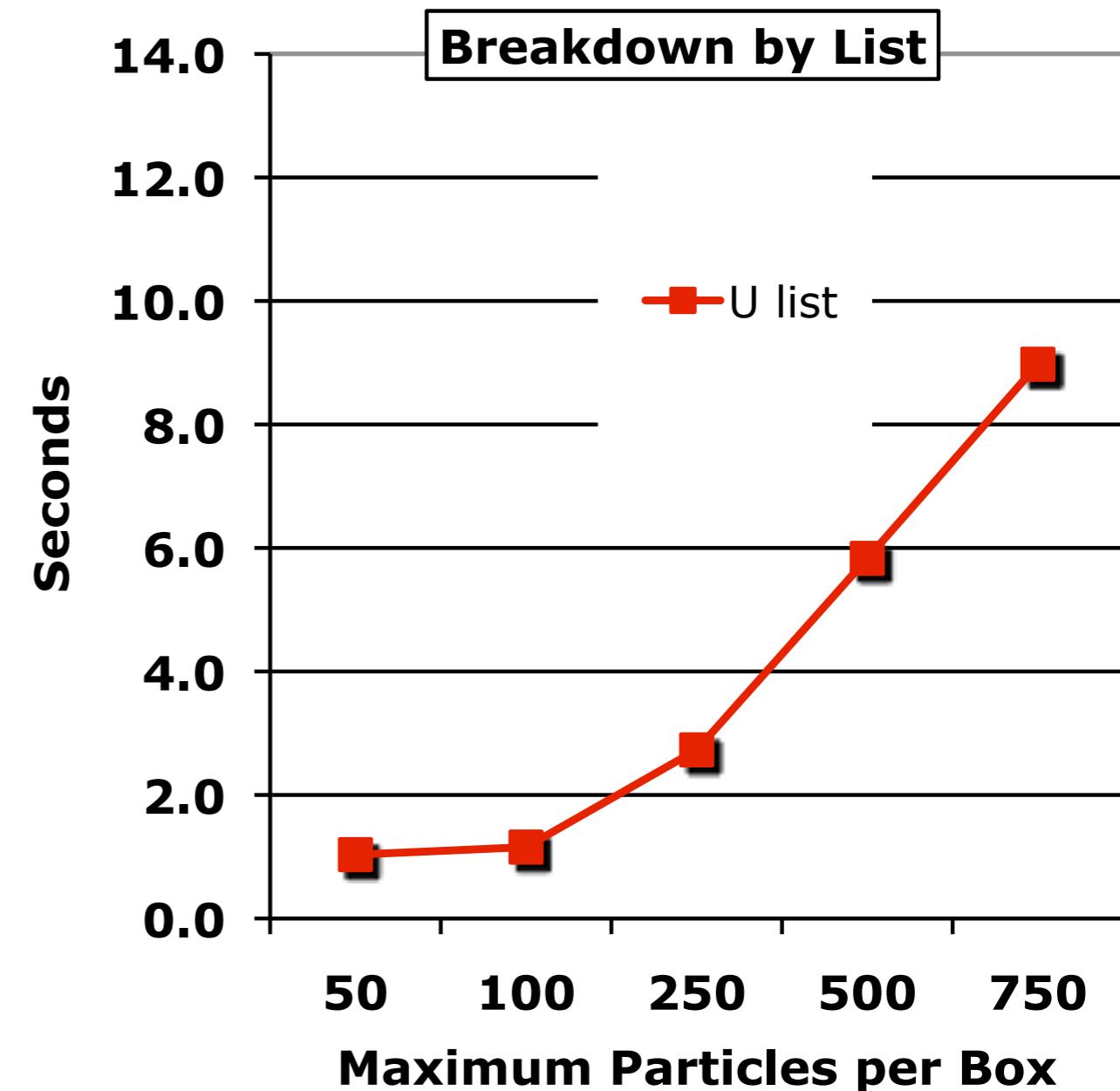
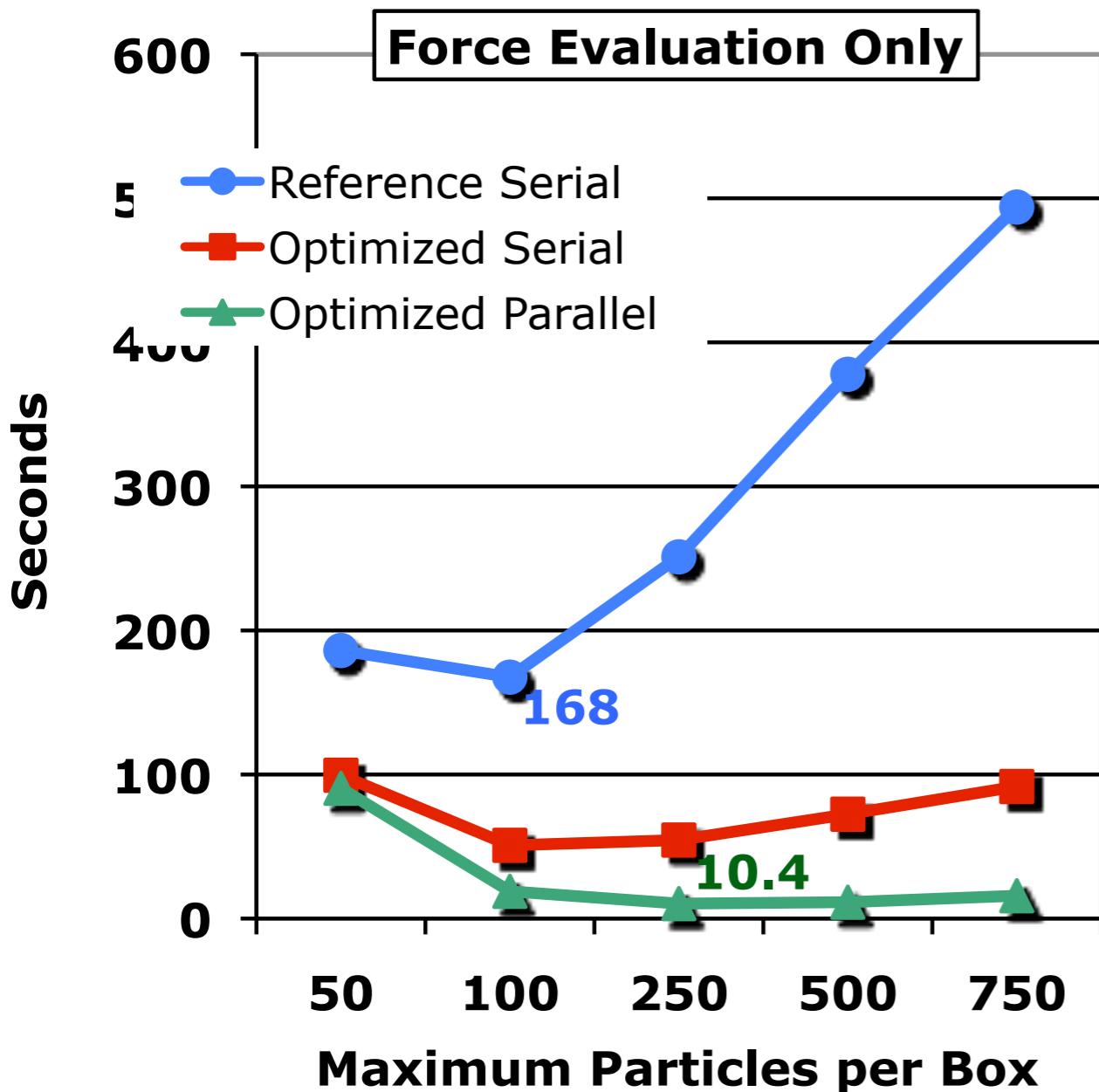
Nehalem



*Shape of curve changes as we introduce optimizations.*

# Algorithmic Tuning of $q = \text{Max pts} / \text{box}$

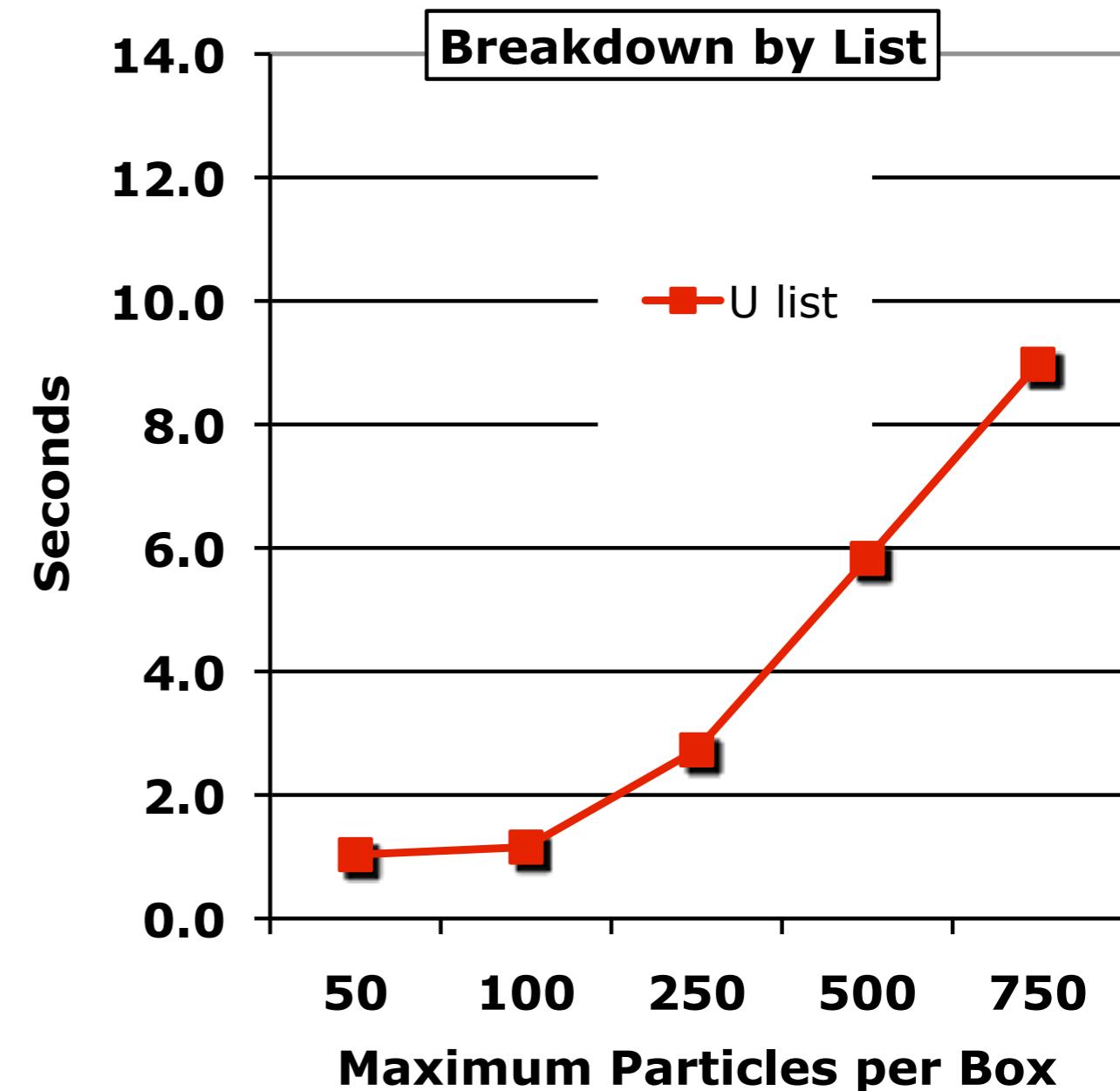
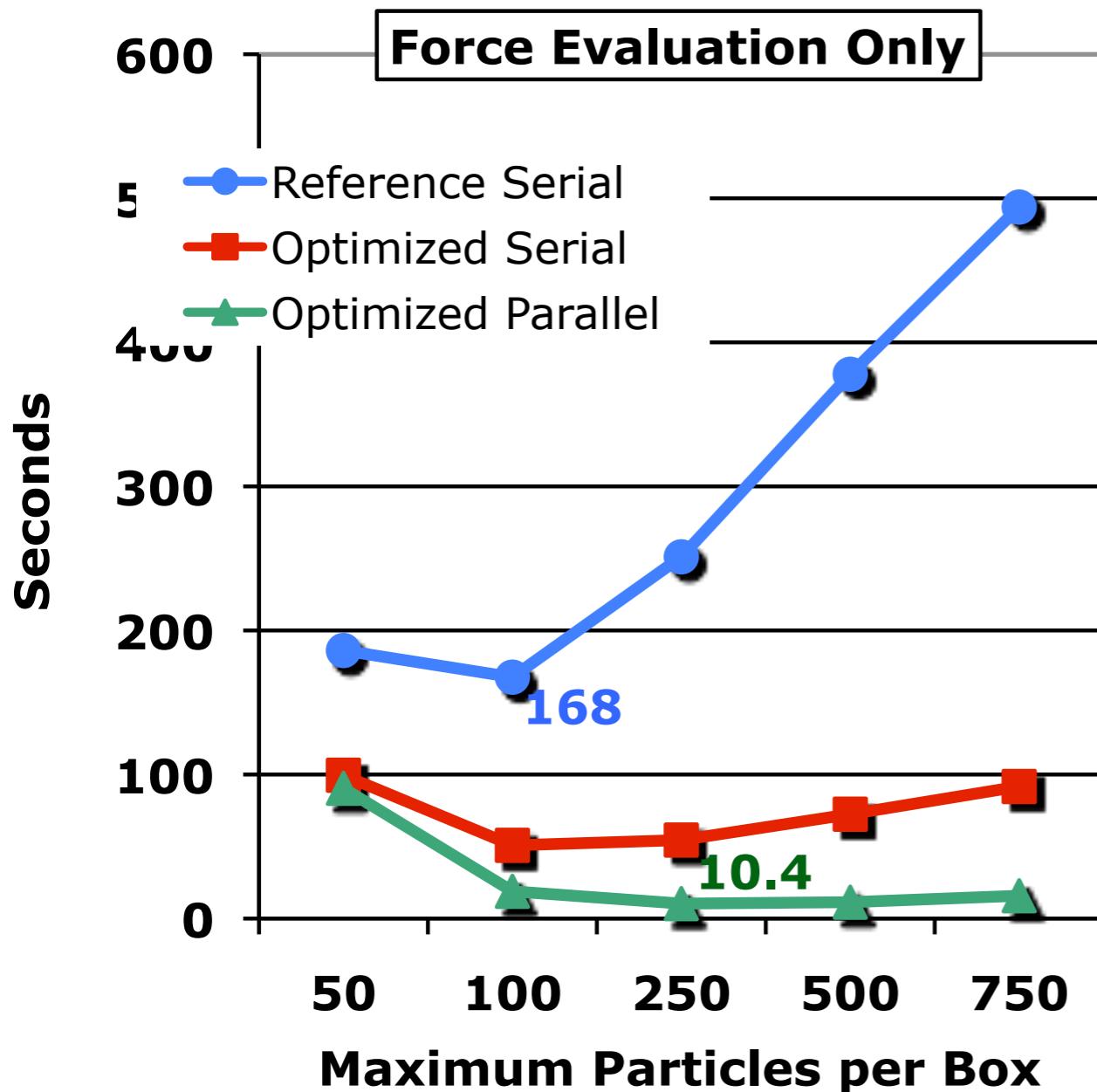
Nehalem



Why? Consider phase costs for the “Optimized Parallel” implementation.

# Algorithmic Tuning of $q = \text{Max pts / box}$

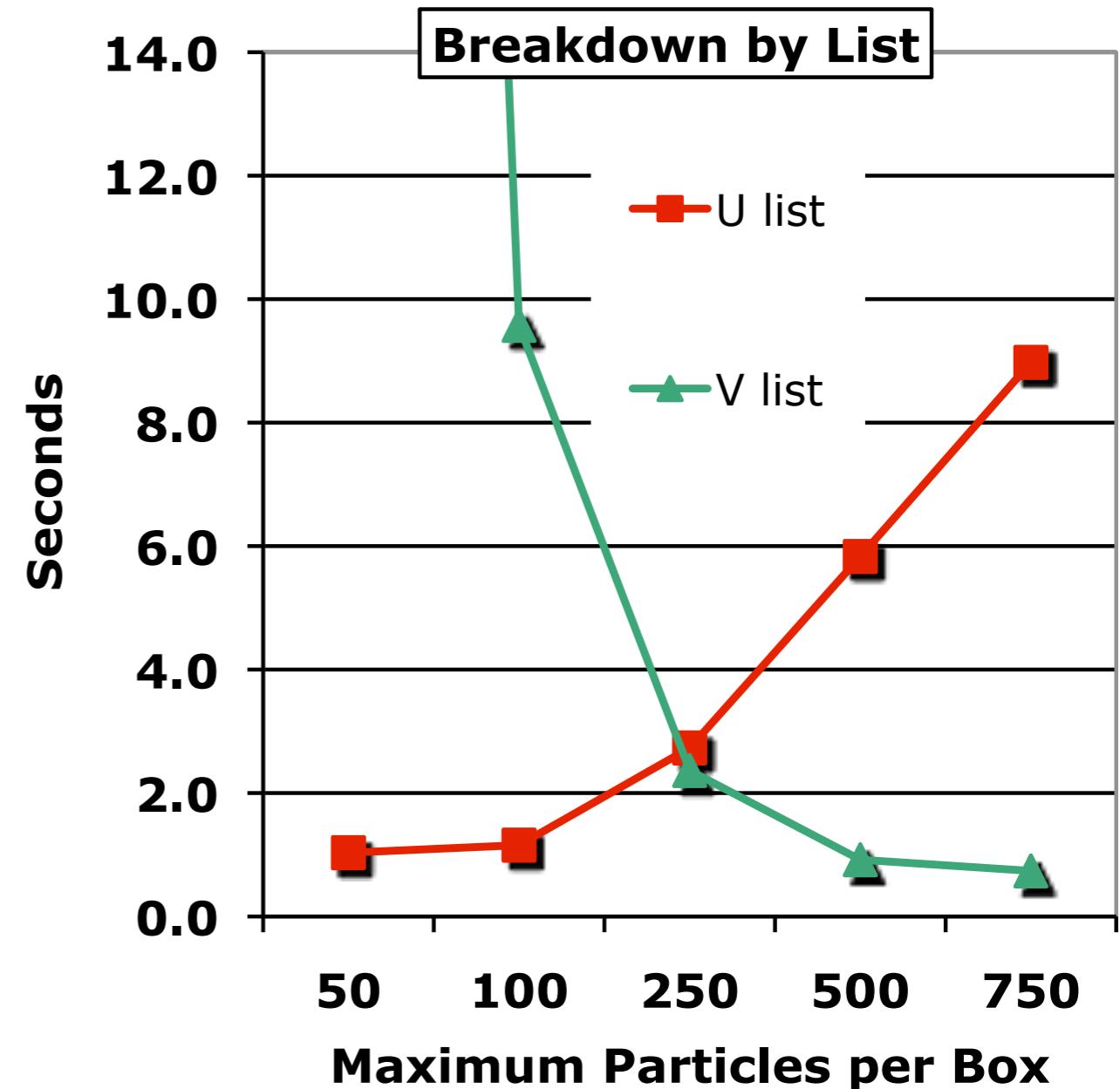
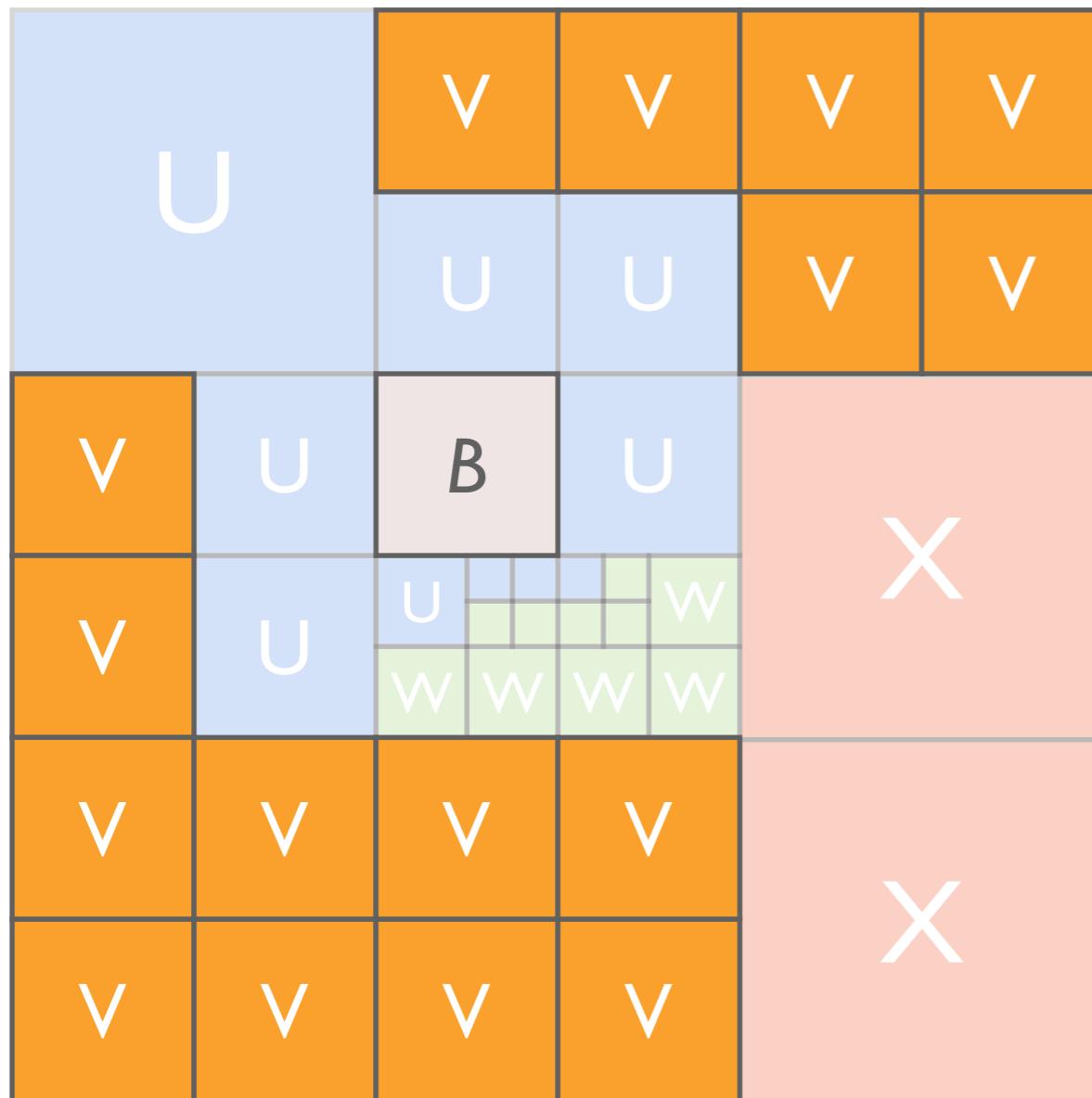
Nehalem



Recall:  $\text{Cost}(U\text{-list}) \sim O(q^2) \text{ per box}$

# Algorithmic Tuning of $q = \text{Max pts / box}$

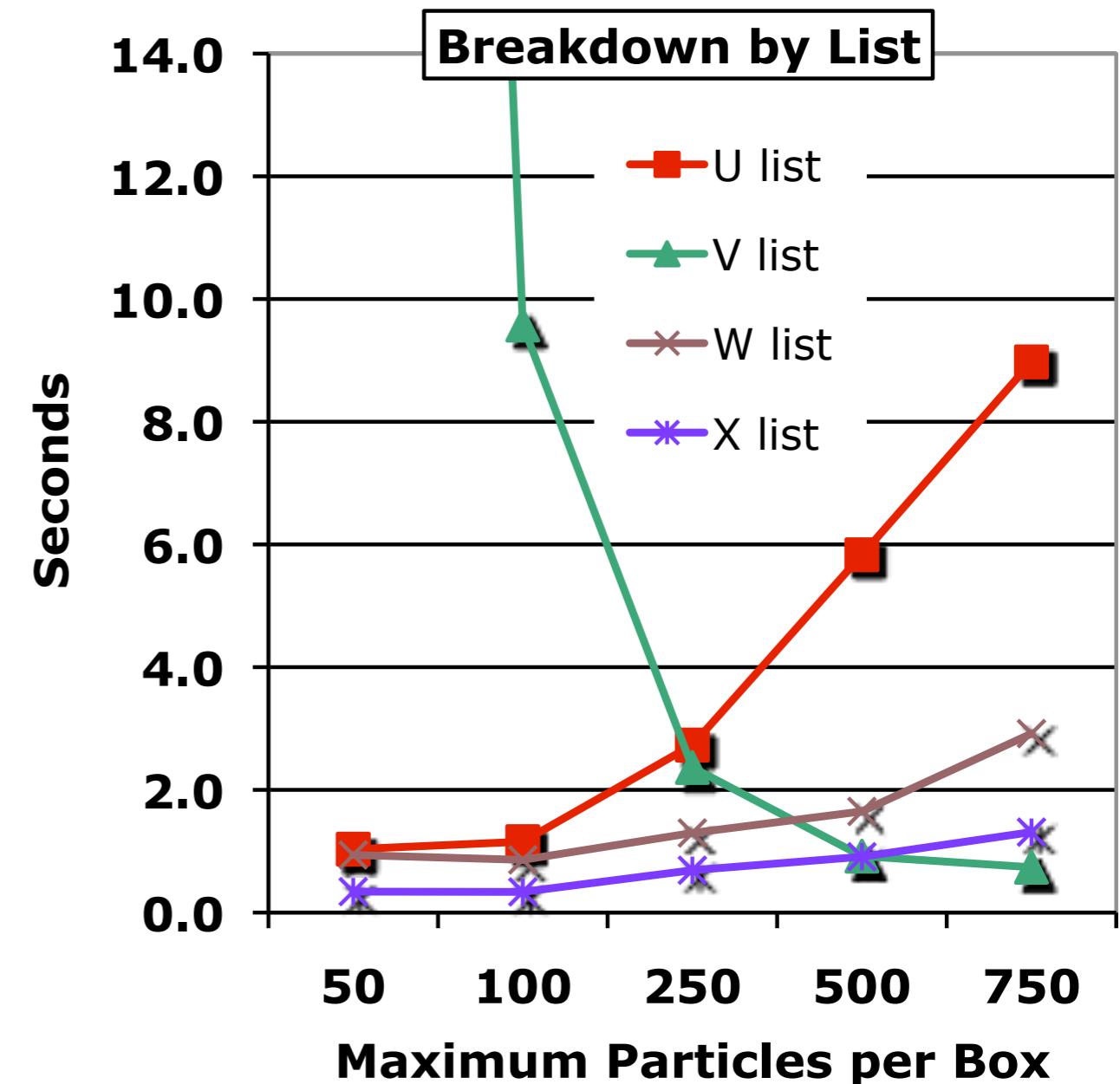
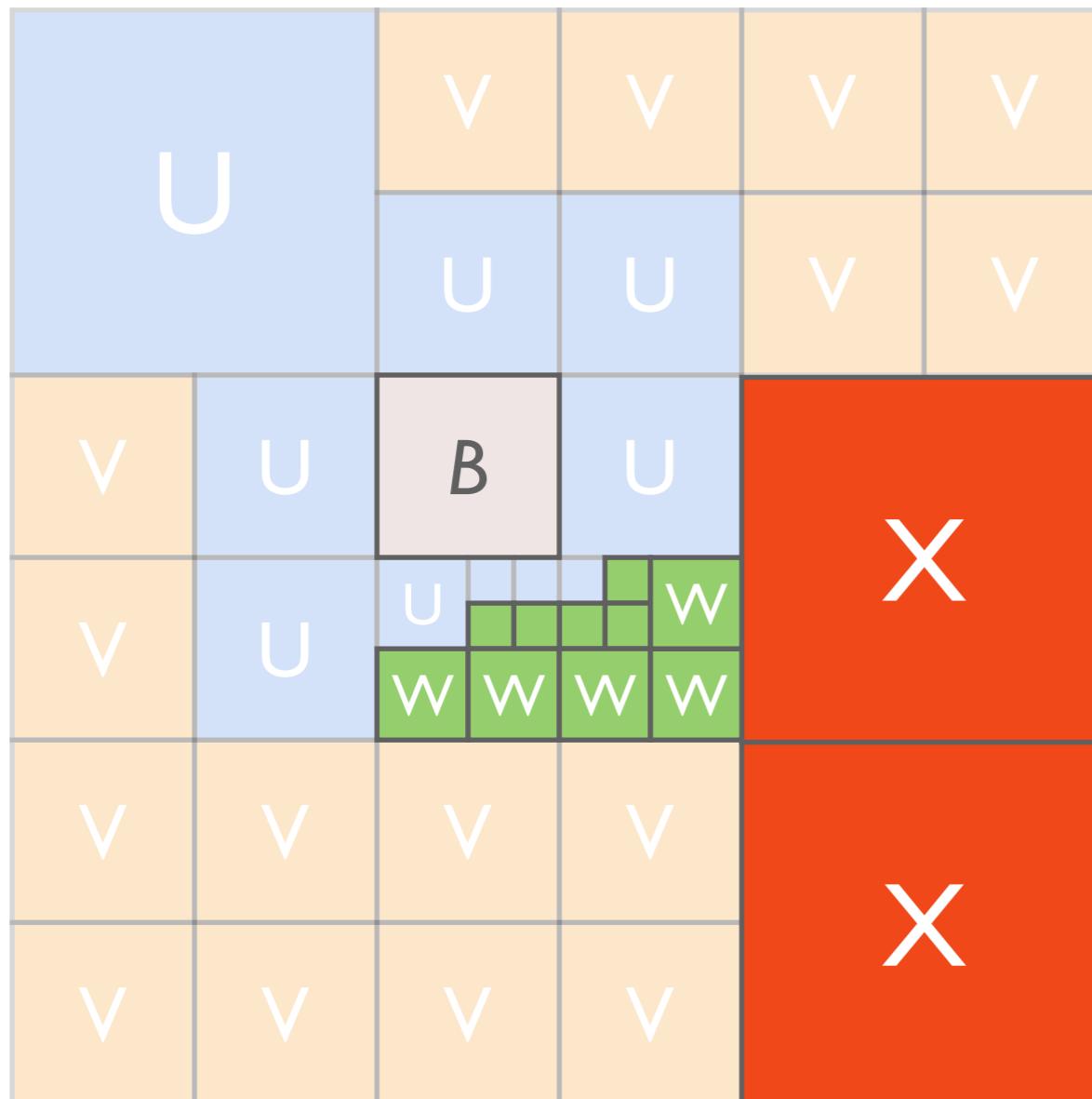
Nehalem



A more shallow tree reduces cost of V-list phase.

# Algorithmic Tuning of $q = \text{Max pts / box}$

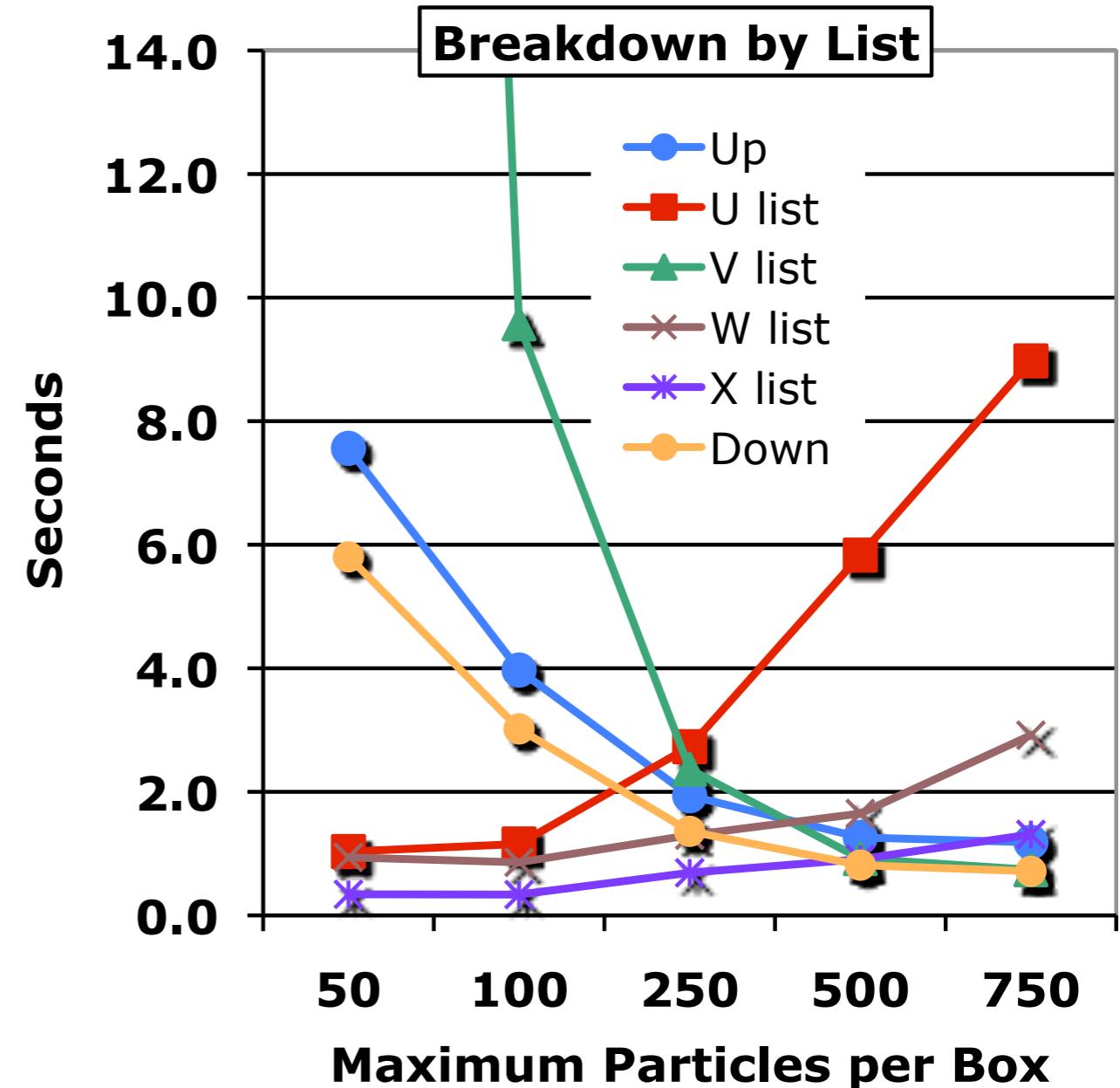
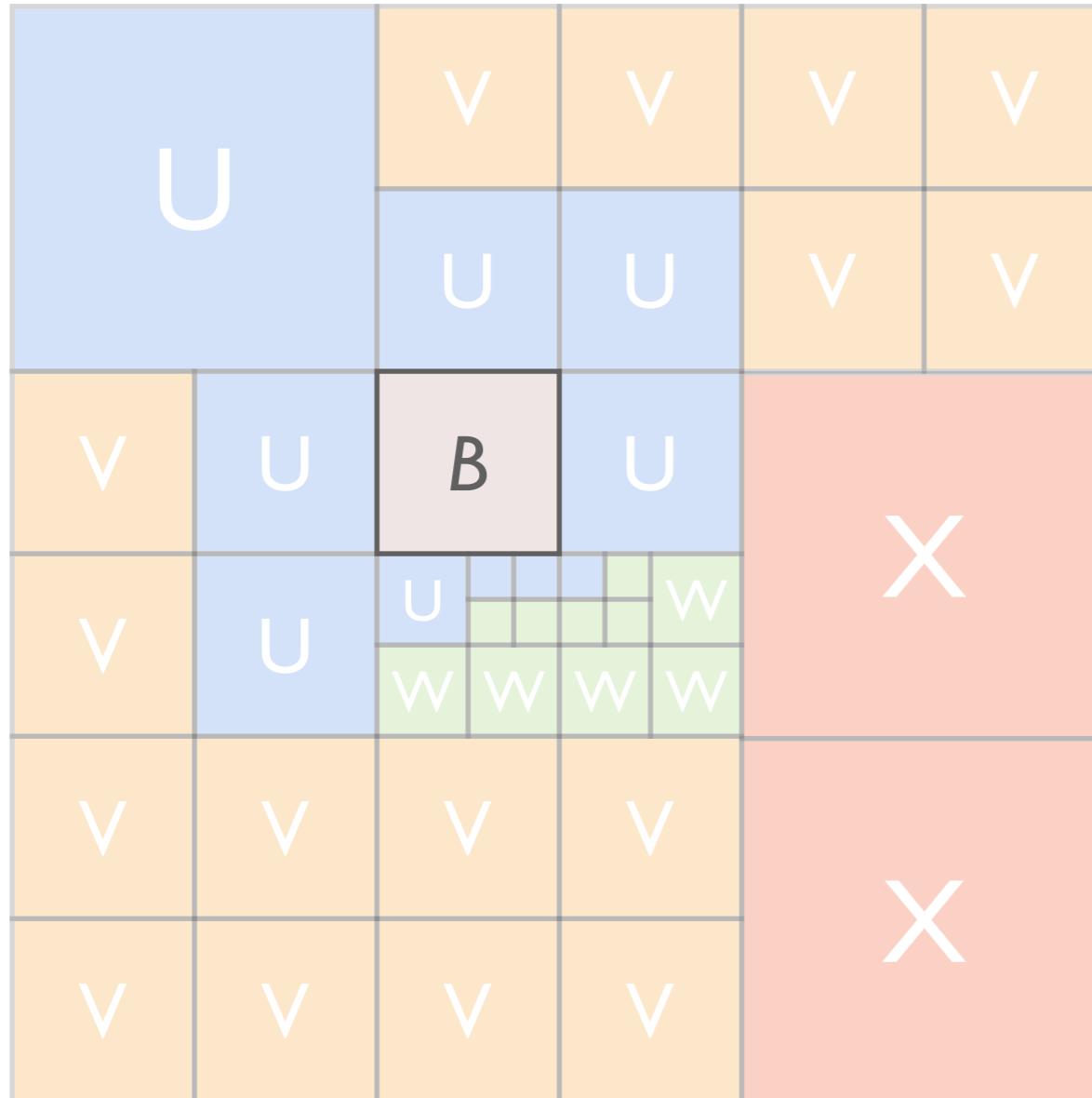
Nehalem



*Computational intensity of W, X more like U than V.*

# Algorithmic Tuning of $q = \text{Max pts / box}$

Nehalem

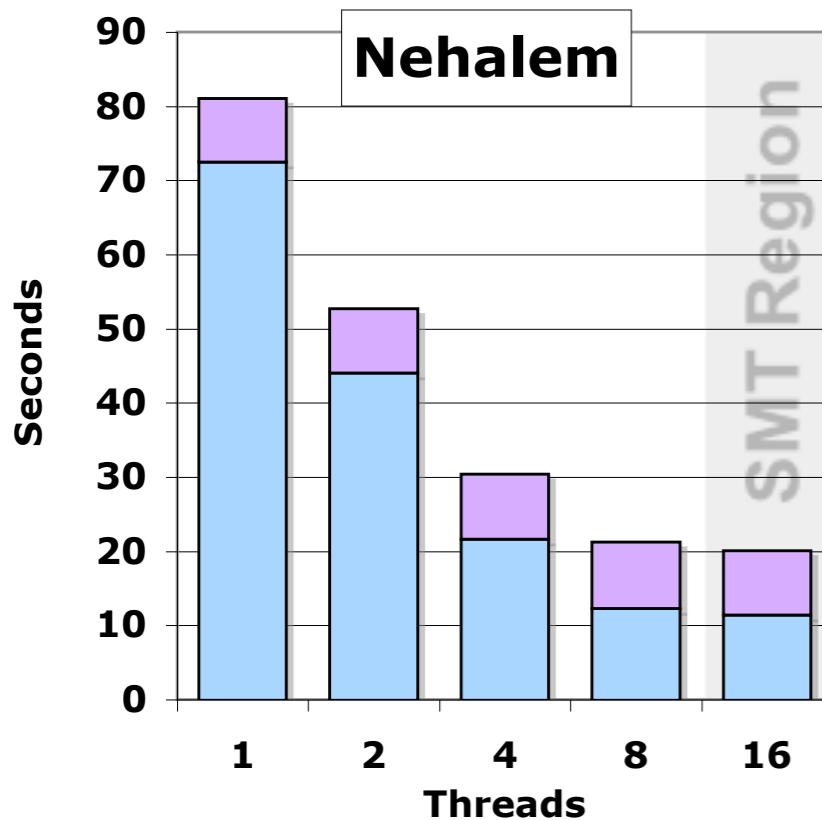


*Optimal q will vary as the point distribution varies.*

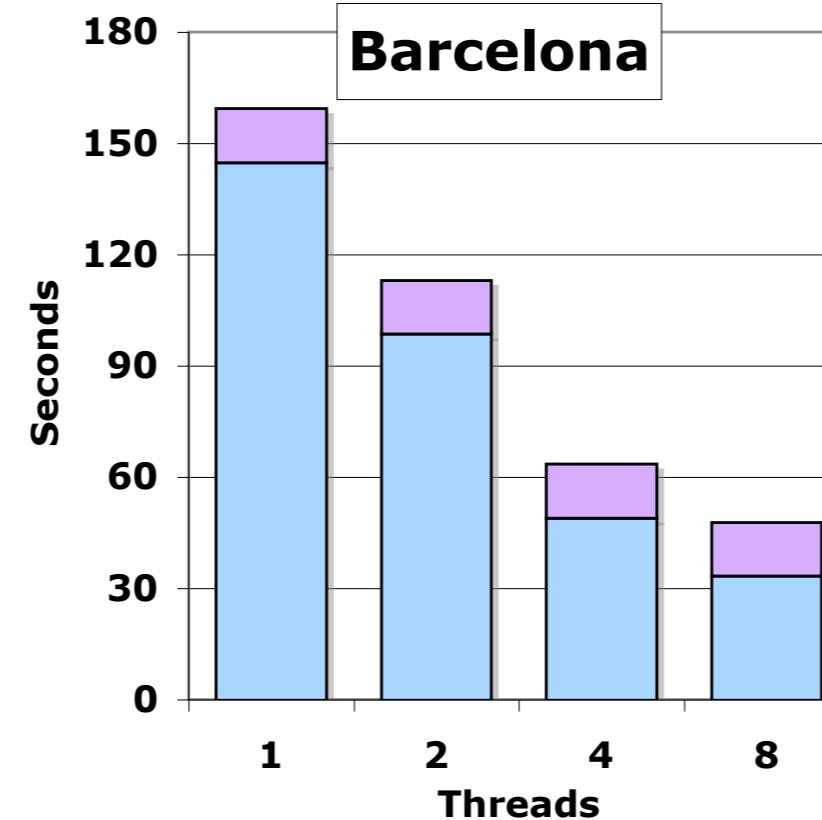
# Multicore Scalability over Optimized Baseline

## Ellipsoidal Distribution

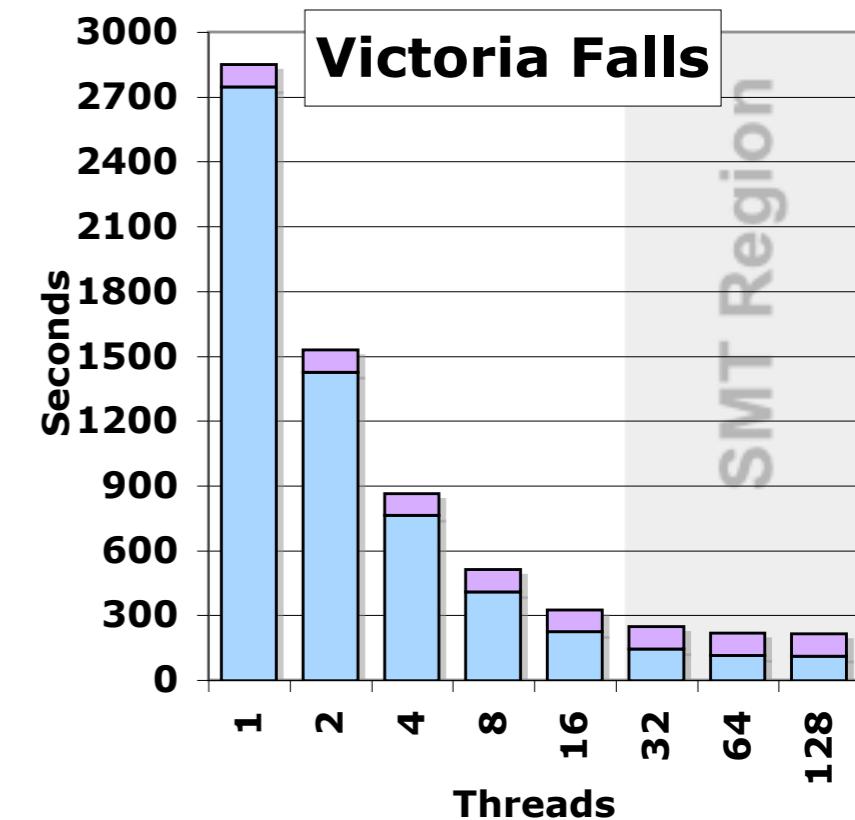
~ 6.3x



~ 4.3x



~ 24x



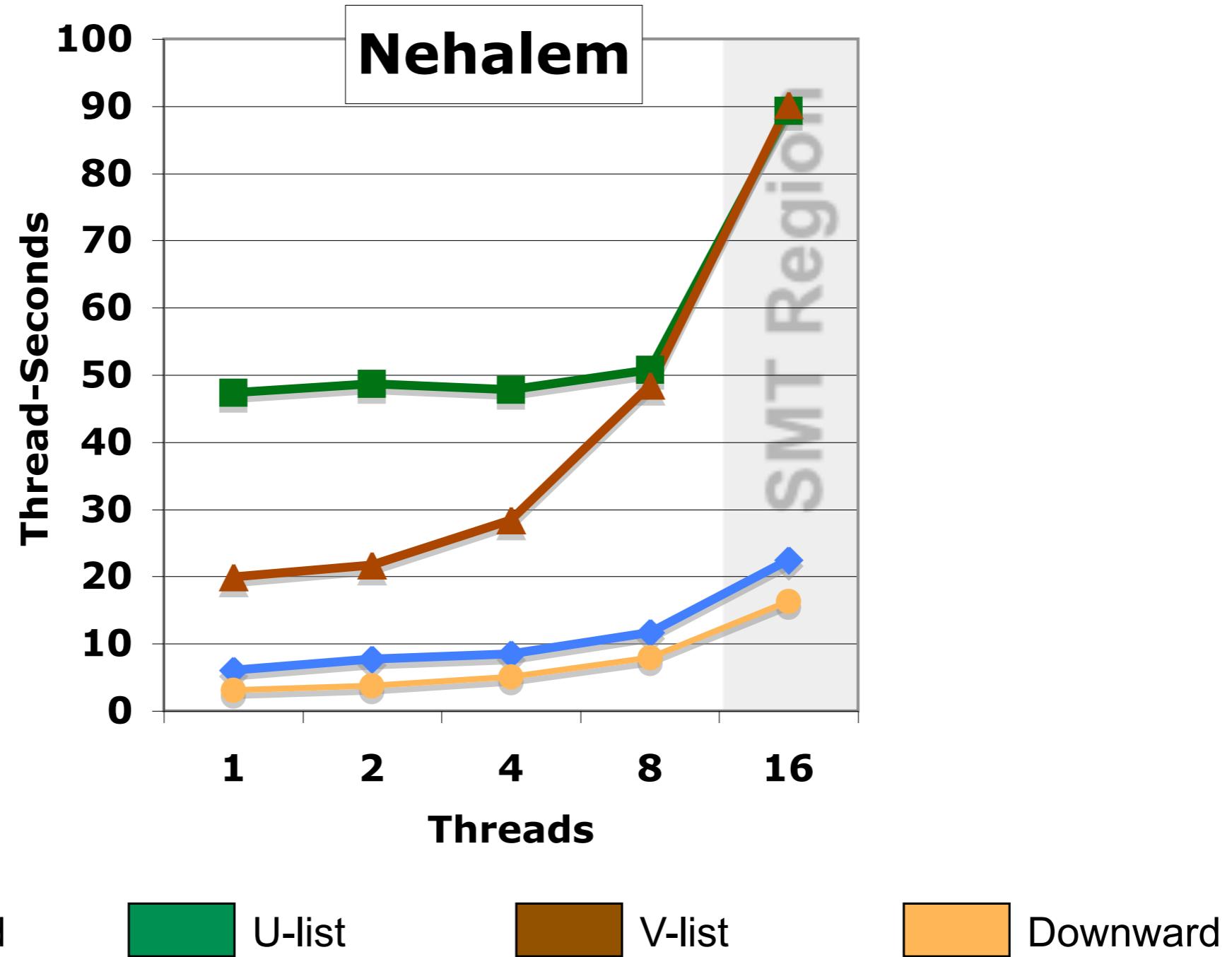
New tree constructed for every force evaluation

asymptotic limit (force evaluation time only)

Need to improve tree construction. Little benefit from SMT.

# Efficiency, via Parallel Cost – $p \cdot T_p$

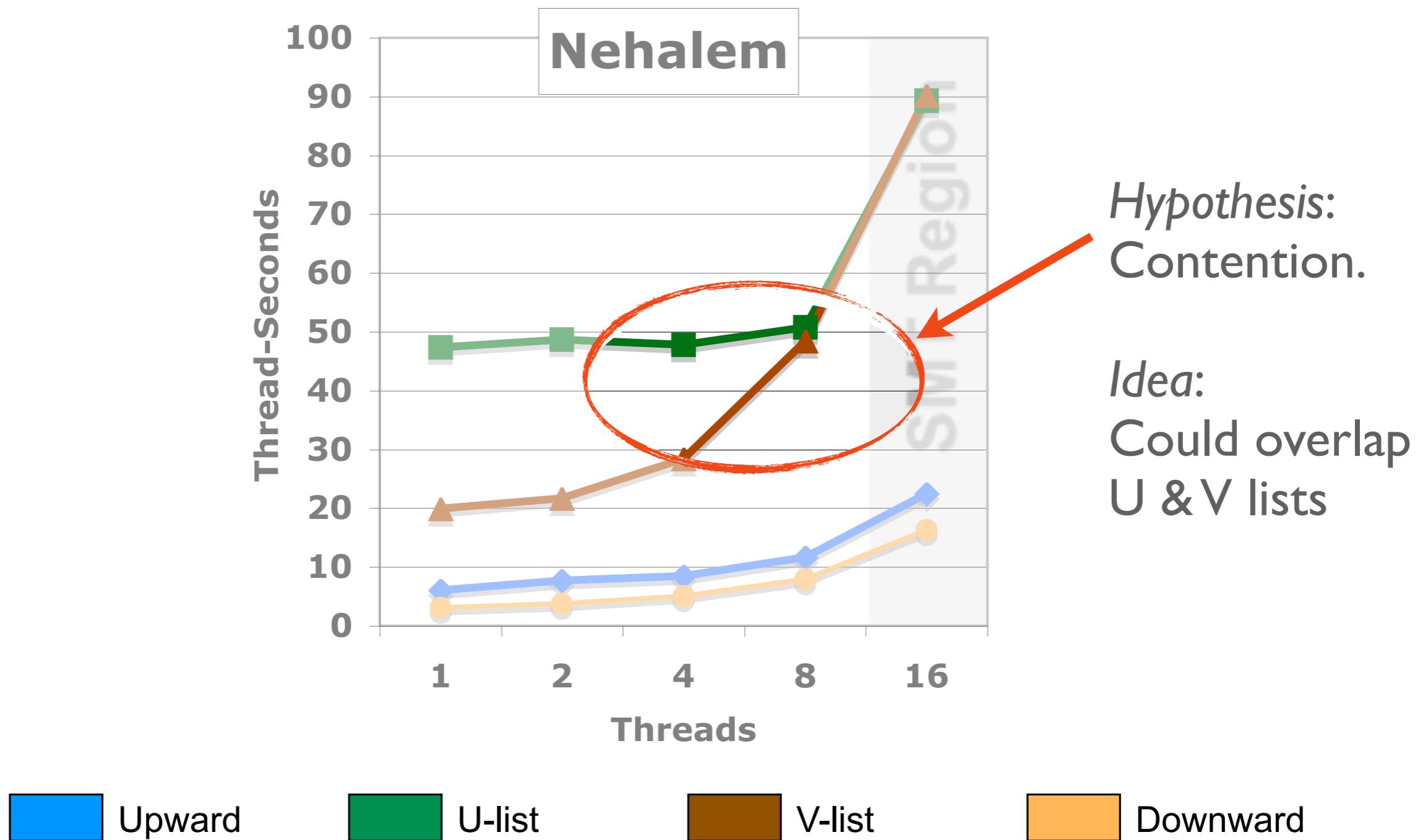
## Uniform Distribution



*Flat horizontal line = perfect scaling*

# Efficiency, via Parallel Cost – $p \cdot T_p$

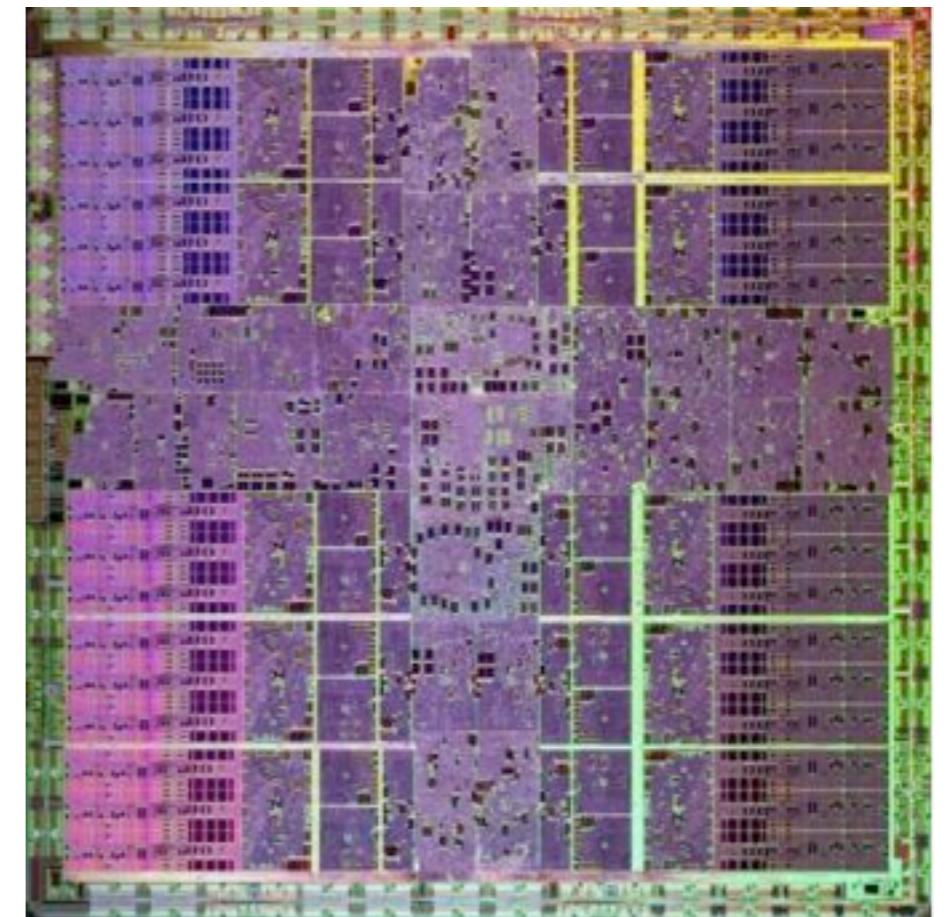
Uniform Distribution



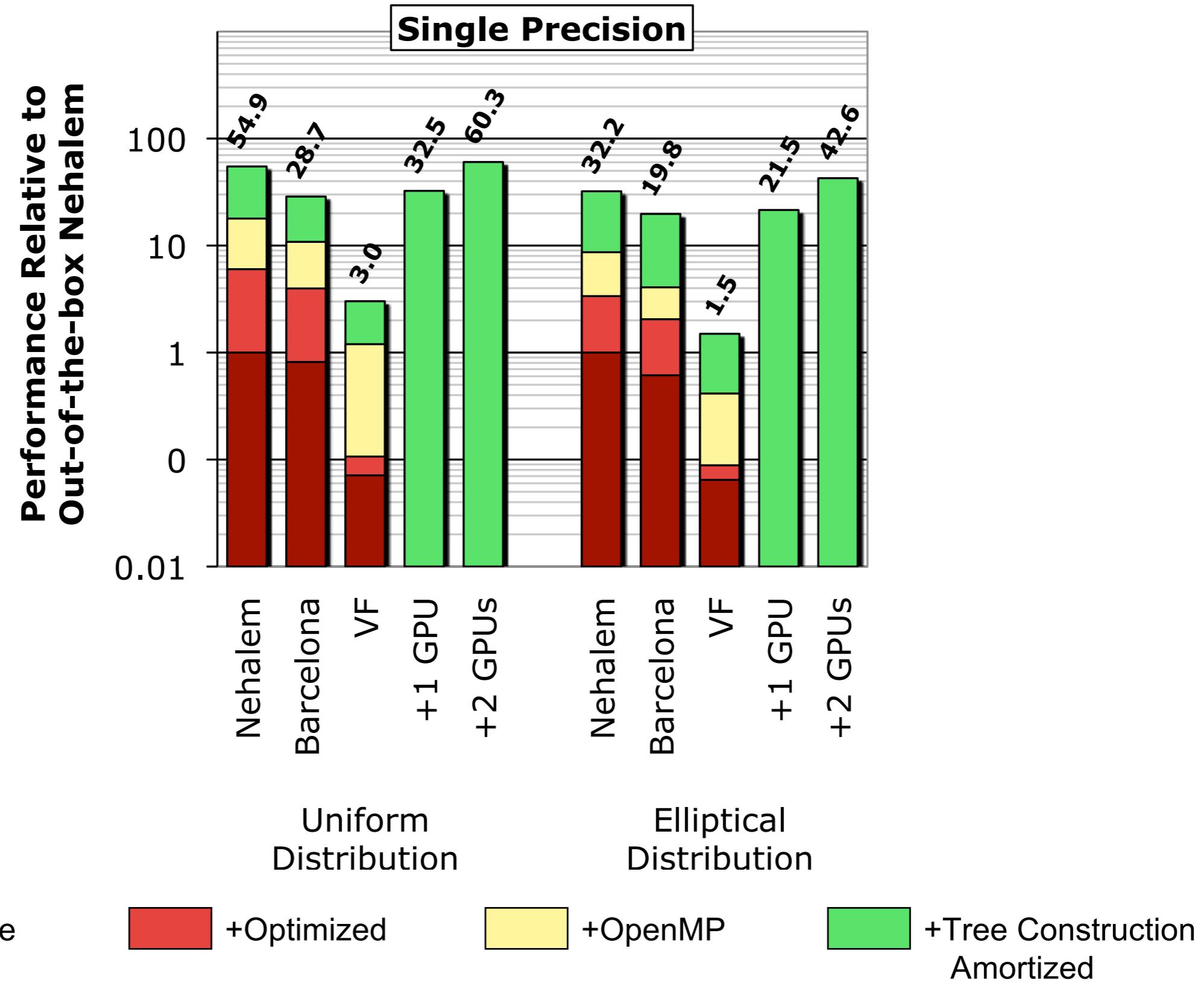
Flat horizontal line = perfect scaling

# GPU comparison: NVIDIA T10P

- ▶ Our prior work on MPI+CUDA  
Lashuk, et al., SC'09
- ▶ System: NCSA Lincoln Cluster
  - ▶ Dual-socket Xeon
  - ▶ 1 node, 1 MPI task per socket & GPU  
(tasks mostly idle)
  - ▶ 1- and 2-GPU configs
  - ▶ Single-precision only for now
- ▶ **12x compute + 5x bandwidth**

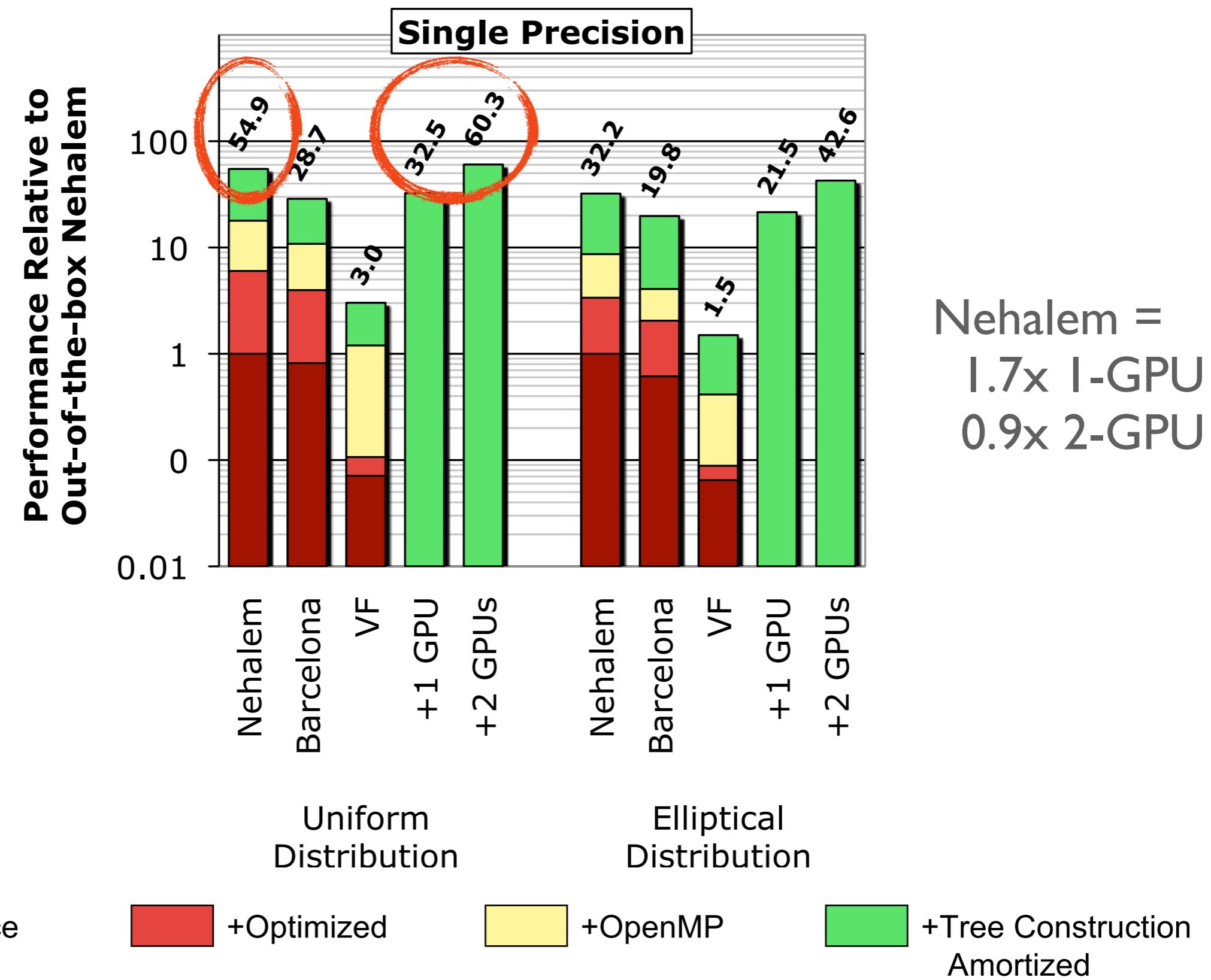


# Cross-Platform Performance Comparison (Summary)



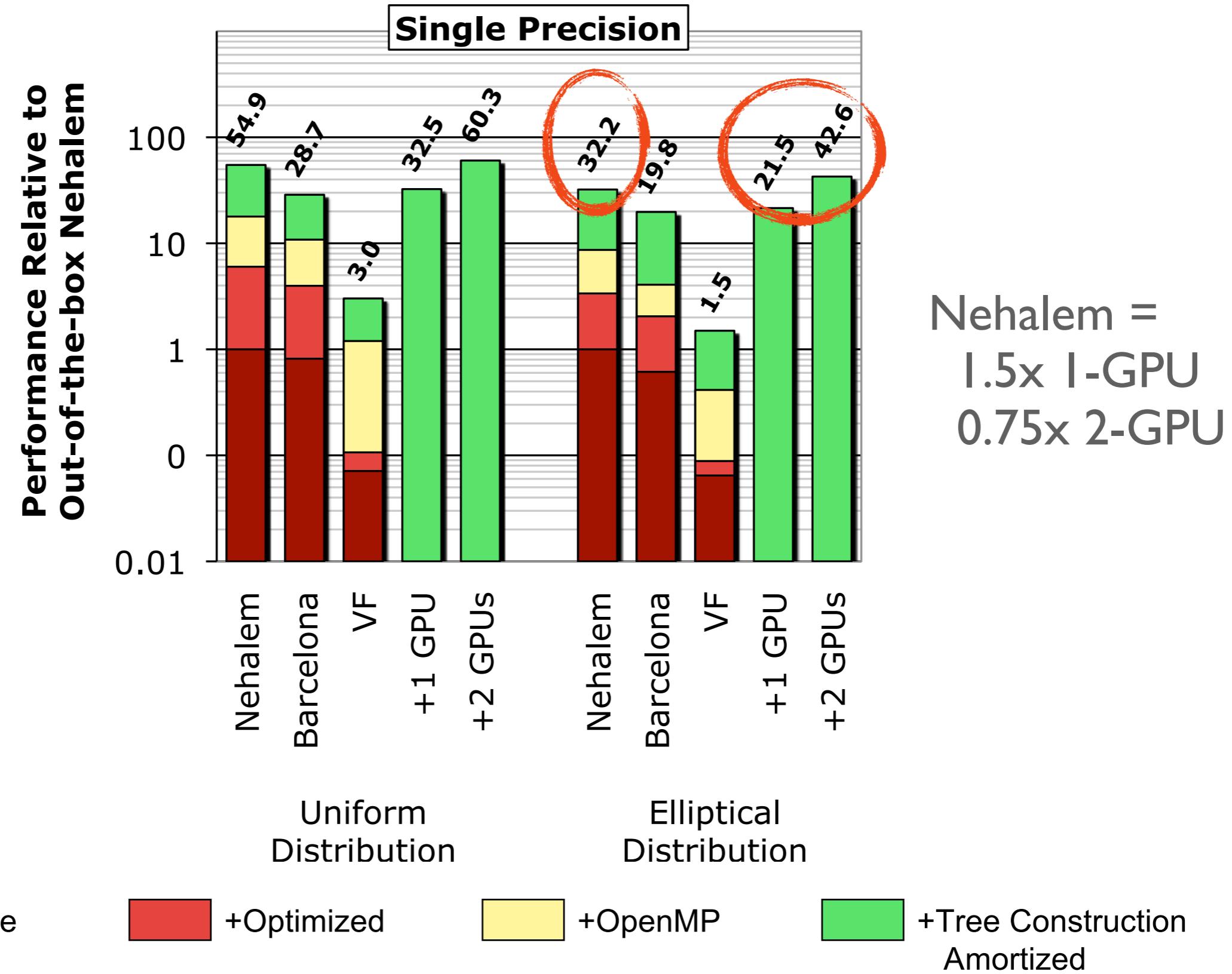
*Nehalem outperforms 1-GPU case, a little slower than 2-GPU case.*

# Cross-Platform Performance Comparison (Summary)



*Nehalem outperforms 1-GPU case, a little slower than 2-GPU case.*

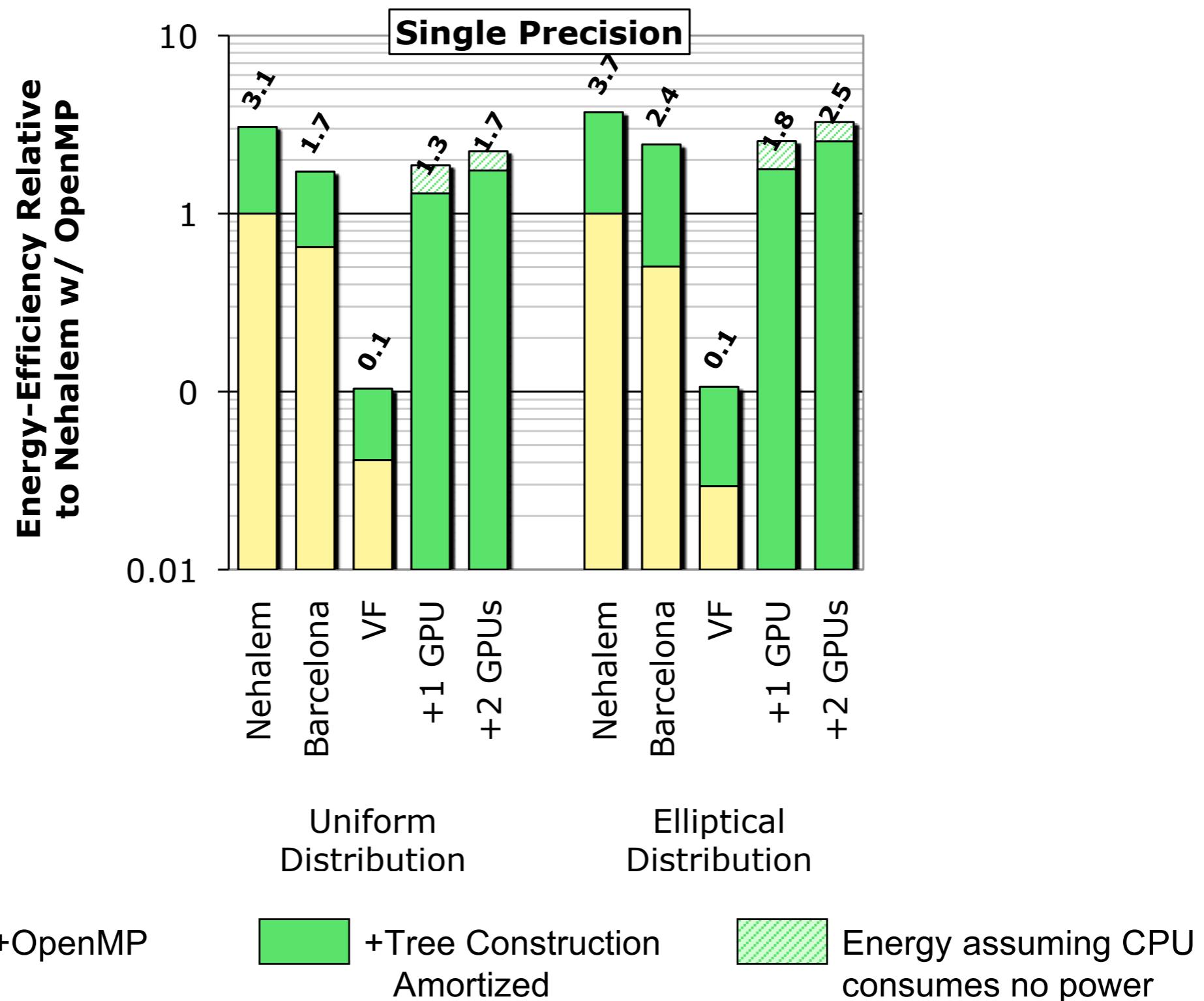
# Cross-Platform Performance Comparison (Summary)



*Nehalem outperforms 1-GPU case, a little slower than 2-GPU case.*

# Cross-Platform Energy-Efficiency Comparison

(Watt-Hours) / (Nehalem+OpenMP Watt-Hours)



Nehalem has same or better power efficiency than either GPU setup.

# Summary and Status

- ▶ First extensive multicore platform study for FMM
  - ▶ Show 25x Nehalem, 9.4x Barcelona, 37.6x VF from algorithmic, data, and numerical tuning
  - ▶ Multicore CPU  $\approx$  GPU in power-performance
- ▶ Short-term:
  - ▶ Perform more detailed modeling → autotuning
  - ▶ Build integrated MPI+CPU+GPU implementation
  - ▶ Parallel tree construction
- ▶ Long-term: Generalize infrastructure and merge with on-going THOR effort for data analysis